

# Class Notes for Biostat 250AB

Academic year 2020-2021: the era of COVID-19 pandemic

Scribe: Elvis Cui, Instructor: Weng Kee Wong

April 20, 2022

## Contents

<b>1</b>	<b>Useful Results in Matrix Theory</b>	<b>5</b>
1.1	Diagonalization and $g$ -Inverse	5
1.2	Partitioning	6
1.3	More on $g$ -inverses	8
1.4	Extreme values of quadratic forms	9
1.5	Projection operator	11
1.6	Exercises	13
<b>2</b>	<b>Distribution Theory</b>	<b>17</b>
2.1	Multivariate normal distribution	17
2.2	Non-central distribution	18
2.3	Basic results for quadratic forms	21
2.4	General results for quadratic forms	23
2.5	Fisher-Cochran's theorem	25
2.6	Applications of Fisher-Cochran's theorem	28
2.7	Exercises	30
<b>3</b>	<b>Theory of Least Square Estimation</b>	<b>33</b>
3.1	Basics of linear models	33
3.2	Properties of LS and G-LS estimators	33
3.3	Adding regressors	34
3.4	More on projections	35
3.5	Constrained least square estimation	38
3.6	When $X$ has less than full column rank	39
3.7	Some testing problems	41
3.8	Exercises	44

<b>4</b>	<b>Multiple and Partial Correlation Coefficient</b>	<b>48</b>
4.1	Multiple correlation coefficient (MCC)	48
4.2	Geometry of LSE	49
4.3	Partial correlation coefficient (PCC)	50
4.4	Exercises	51
<b>5</b>	<b>Violation of Assumptions and Remedies</b>	<b>52</b>
5.1	Outlier detection	52
5.2	Under-fitting and over-fitting	52
5.3	Transformation	54
5.4	Collinearity	55
5.5	Exercises	56
<b>6</b>	<b>Hypothesis Testing and Simultaneous Inference</b>	<b>58</b>
6.1	Fieller's theorem	58
6.2	Lack of fitness test	61
6.3	Scheffe's method	62
6.4	Tukey's $q$	63
6.5	Test of homogeneity	65
6.6	Cook-Weisberg score test	66
6.7	Exercises	68
<b>7</b>	<b>Shrinkage and Bayes Estimation</b>	<b>70</b>
7.1	Principal component regression	70
7.2	Ridge estimator	70
7.3	James-Stein estimator	72
7.4	Bayes estimator	73
7.5	Exercises	75
<b>8</b>	<b>ANOVA Mixed Models</b>	<b>77</b>
8.1	One-way ANOVA with random effects	77
8.2	Two factor ANOVA	79
8.3	Satterwaite approximation	80
8.4	Exercises	81
<b>9</b>	<b>Linear Mixed Models</b>	<b>82</b>
9.1	Henderson's equation	82
9.2	Best linear unbiased predictor (BLUP)	83
9.3	Restricted maximum likelihood (REML)	86
9.4	Exercises	88

<b>10 Spline Regression</b>	<b>89</b>
10.1 Regression splines . . . . .	89
10.2 B-splines . . . . .	91
10.3 Smoothing splines . . . . .	95
10.4 The selection of knots position . . . . .	96
10.4.1 Conventional methods . . . . .	96
10.4.2 Bayesian methods . . . . .	96
10.4.3 Metaheuristic methods . . . . .	96

# Preface

This set of lecture notes is based on the 2 courses I took from 2020 to 2021 (Biostat 250A and 250B). From 2020 to 2021, all my classmates (and colleagues under the assumption that TA is an official employee of UCLA) took and taught classes online, which is, I believe, the **VERY FIRST TIME** in history ever. COVID-19 pandemic has changed a lot to our daily life, not only teaching and taking classes, but also working, communicating with other, etc. What stays unchanged is our passion to learn and study.

Note that  $x$  and  $\mathbf{x}$  can both represent a vector and  $X$  and  $\mathbf{X}$  can both represent a matrix. Also, I switch between  $y$  and  $Y$  from time to time. Both  $\mathcal{N}$  and  $\mathcal{N}_p$  can represent a multivariate distribution. I hope they do not cause too much confusion and I am too lazy to modify them so I apologize. References are given and some key results such as generalized Cauchy-Schwarz inequality and Cook-Weisberg score test are based them.

Finally, a quote from CR Rao,

$$\text{Statistics} = \text{Approximation} + \text{Optimization}$$

and another quote from George Box,

All models are wrong, some are useful.

Elvis Cui

July 5, 2021.

# 1 Useful Results in Matrix Theory

This section is not based on any single lecture but a collection of all results in matrix theory that I learned in Biostat 250AB and other related materials.

## 1.1 Diagonalization and g-Inverse

**Theorem 1.1** (Simultaneous diagonalization theorem). Let  $A, M$  be two symmetric  $n \times n$  matrices,  $M$  positive definite. Then  $\exists C$  and  $\det |C| \neq 0$  s.t.

$$C^T M C = I, C^T A C = \Lambda$$

where  $I$  is identity and  $\Lambda$  is diagonal.

*Proof.*  $M \succ 0 \Rightarrow M = M^{1/2} M^{1/2}$  and let  $B := (M^{-1/2})^T A (M^{-1/2})$ .

Then  $\exists P^T$   $n' \times n$  orthogonal, i.e.,  $PP^T = I$  s.t.  $PBP^T = \Lambda$  or  $P^T M^{-1/2} A M^{-1/2} P = \Lambda$ . The diagonal elements of  $\Lambda$  are eigenvalues of  $B$ . And let  $C = M^{-1/2} P$ , then we have  $\det |C| \neq 0$  and

$$C^T A C = \Lambda, C^T M C = I$$

□

**Theorem 1.2** (Existence of g-inverse). Let  $X$  be a  $n \times p$  matrix and  $\text{rank}(X) = r < \min(n, p)$ , then there exists infinitely many  $X^-$ , they are  $p \times n$  matrices and satisfy

$$X X^- X = X$$

*Proof.* Let  $X = S D T^T$  be the singular value decomposition and  $D$  is a  $r \times r$  matrix. Define

$$X^- = T D^{-1} S^T$$

and we have  $X X^- X = S D T^T T D^{-1} S^T S D T^T = S D T^T = X$ . But there are infinitely many g-inverses:

$$\tilde{X} = X^- + (I - X^- X) B$$

where  $B$  is any  $p \times n$  matrix. We have

$$X \tilde{X} X = X X^- X + X (I - X^- X) B X = X + X B X - X X^- X B X = X$$

□

**Example 1.1** (Solution of linear systems). Consider the linear system

$$A x = b$$

A solution  $x^*$  exists iff it is a consistent set of equations, i.e.,

$$\text{rank}(A|b) = \text{rank}(A)$$

and

$$x^* = A^{-}b$$

is a solution.

**Theorem 1.3** (Positive definite relation). If  $A \succ B$  and both  $A, B$  are positive definite, then we have

$$B^{-1} \succ A^{-1}$$

*Proof.* By theorem 1.1, there exists a non-singular matrix  $U$  such that

$$A^{-1} = U^T U, \quad B^{-1} = U^T D^{-1} U$$

where  $D$  is diagonal with elements equal to eigenvalues of  $A^{-1/2} B A^{-1/2}$ . Then,

$$A - B = U^{-1}(I - D)U^{-T} \succ 0$$

suggests that  $d_{ii} \leq 1$  for all diagonal elements of  $D$ . Thus,

$$A^{-1} - B^{-1} = U^T(I - D^{-1})U \prec 0$$

since  $(I - D^{-1}) \prec 0$ .

□

## 1.2 Partitioning

**Theorem 1.4** (Determinant formula). Let

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}, \quad B = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$$

Then:

1. If  $A_{21}=0$ , then

$$\left| \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \right| = |A_{11}| |A_{22}|$$

2. If  $A$  and  $B$  are square matrices of the same dimension, then

$$|I + AB| = |I + BA|$$

3. In the case  $A_{21} \neq 0$ , we have

$$\left| \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \right| = |A_{22}| |A_{11} - A_{12}A_{22}^{-1}A_{21}|$$

*Proof.* I only prove 2.

$$\begin{pmatrix} I & A \\ -B & I \end{pmatrix} \begin{pmatrix} I & 0 \\ B & I \end{pmatrix} = \begin{pmatrix} I + AB & A \\ 0 & I \end{pmatrix}$$

$$\begin{pmatrix} I & 0 \\ B & I \end{pmatrix} \begin{pmatrix} I & A \\ -B & I \end{pmatrix} = \begin{pmatrix} I & A \\ 0 & I + BA \end{pmatrix}$$

But the determinants of LHS are identical.

□

**Theorem 1.5** (Inverse of block matrices). Let

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

Then

$$B_{11} = A_{11.2}^{-1} = (A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1}$$

$$B_{21} = -A_{22}^{-1}A_{21}A_{11.2}^{-1}$$

$$B_{12} = -A_{11.2}^{-1}A_{12}A_{22}^{-1}$$

$$B_{22} = A_{22}^{-1} + A_{22}^{-1}A_{21}A_{11.2}^{-1}A_{12}A_{22}^{-1}$$

**Theorem 1.6** (Diagonal elements of leverage matrix, [1]). Let  $W$  be an  $p \times p$  matrix of rank  $p$  with columns  $w_j, j = 1, \dots, p$ . Then

$$(W^T W)_{jj}^{-1} = [w_j^T (I - P_{-j}) w_j]^{-1} \quad (1.1)$$

where  $P_{-j}$  is the projection operator associated with  $W_{-j}$ , i.e.,  $W$  with its  $j^{th}$  column omitted.

*Proof.* By inverse of block matrices 1.5, we have

$$\begin{aligned}(W^T W)^{-1} &= \begin{pmatrix} W_{-p}^T W_{-p} & W_{-p}^T w_p \\ w_p^T W_{-p} & w_p^T w_p \end{pmatrix}^{-1} \\ &= \begin{pmatrix} F & g \\ g^T & h \end{pmatrix}\end{aligned}$$

where  $h = (W^T W)_{pp}^{-1} = (w_p^T w_p - w_p^T P_{-p} w_p)^{-1}$ . Let  $\Pi$  be the permutation matrix  $I$  with its  $j^{th}$  and  $p^{th}$  columns interchanged. Then  $\Pi^2 = I$ , so  $\Pi$  is a symmetric orthogonal matrix, and its own inverse. Hence

$$\begin{aligned}(W^T W)^{-1} &= \Pi(\Pi W^T W \Pi)^{-1} \Pi \\ &= \Pi \begin{pmatrix} F_1 & g_1 \\ g_1^T & h_1 \end{pmatrix} \Pi\end{aligned}$$

where  $h_1 = (W^T W)_{jj}^{-1}$ . Thus  $w_p$  and  $w_j$ ,  $W_{-p}$  and  $W_{-j}$  have been effectively interchanged. The result then follows.  $\square$

### 1.3 More on $g$ -inverses

Moore (1935) and Penrose defined a special type of generalized inverse  $A^+$  known as Moore-Penrose pseudo-inverse. It satisfies the following 4 properties [2].

1.  $AA^+A = A$ . This is the definition of the usual  $g$ -inverse.
2.  $A^+AA^+ = A^+$ .
3.  $(AA^+)^T = AA^+$ .
4.  $(A^+A)^T = A^+A$ .

Such inverse exists and is unique (exercise).

Next, suppose  $A$  is a  $p \times p$  block matrix of the following form:

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$$

and  $\text{rank}(A) = \text{rank}(A_{11}) = r < p$ . By **Analissa's transformation** in Biostat 250C, we have

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

But the rank of LHS is the same as  $\text{rank}(A) = r$ . Therefore, the bottom-right block matrix of RHS



must be 0, that is,

$$A_{22} = A_{21}A_{11}^{-1}A_{12}$$

By this identity, we can verify that

$$A^- = \begin{pmatrix} A_{11}^{-1} & O \\ O & O \end{pmatrix} \quad (1.2)$$

is indeed a  $g$ -inverse of  $A$ , i.e.,  $AA^-A = A$ . Further, we also have  $A^-AA^- = A^-$ . Some other useful results on  $g$ -inverse are listed below for further reference, most of them can be verified easily.

**Lemma 1** (Properties of  $g$ -inverse). *For any  $g$ -inverse, the following holds:*

1.  $A(A^TA)^-A^TA = A$  and  $A^TA(A^TA)^-A = A^T$ . We can apply this result to show the estimability of a linear function of  $\beta$  in a classical linear model in chapter 3.
2.  $BA^-B = B$  holds if and only if  $\mathcal{C}(B^T) = \mathcal{C}(A^T)$ .
3. For any choice of  $g$ -inverse  $(A^TA)^-$ , the orthogonal projection  $P = A(A^TA)^-A^T$  is symmetric and does not depend on  $(A^TA)^-$ . We talk more about orthogonal projection matrices in next few sections and in chapter 3.

Next, if  $A = A^T$  is symmetric, then do we have  $(A^-)^T = A^-$ ? By definition,

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

Hence,

$$A^- = (A^T)^- = (A^-)^T. \quad (1.3)$$

## 1.4 Extreme values of quadratic forms

**Theorem 1.7** (A mini-max theorem [3]). Let  $A = A^T$  be a  $n \times n$  symmetric matrix and the spectral decomposition is  $A = P\Lambda P^T$ .  $x$  is a  $n \times 1$  vector and  $B$  is a  $n \times k$  matrix. Then we have

$$\inf_B \sup_{B^Tx=0} \frac{x^T Ax}{x^T x} = \lambda_{k+1}$$

$$\sup_B \inf_{B^Tx=0} \frac{x^T Ax}{x^T x} = \lambda_{n-k}$$

where  $\lambda_i$  is the  $i^{th}$  diagonal element of  $\Lambda$ .

*Proof.* Let  $A = P\Lambda P^T$  be its spectral decomposition and  $P$  has columns  $p_1, p_2, \dots, p_n$ . Then  $x$  admits the representation

$$x = \sum_{i=1}^n c_i p_i = Pc$$

for some  $c^T = (c_1, \dots, c_n)$ . Then  $B^T x = 0$  iff  $B^T P c = 0$  or  $G^T c = 0$  where  $G = P^T B$ . Hence, the problem reduces to

$$\inf_G \sup_{G^T c = 0} \frac{c_1^2 \lambda_1 + \dots + c_n^2 \lambda_n}{c_1^2 + \dots + c_n^2}$$

Note that by letting  $c^T = (c_1, \dots, c_{k+1}, 0, \dots, 0)$ , we have

$$\sup_{G^T c = 0} \frac{c_1^2 \lambda_1 + \dots + c_n^2 \lambda_n}{c_1^2 + \dots + c_n^2} \geq \sup_{G^T c = 0} \frac{c_1^2 \lambda_1 + \dots + c_{k+1}^2 \lambda_{k+1}}{c_1^2 + \dots + c_{k+1}^2} \geq \lambda_{k+1}$$

The equality shall be attained when  $c_{k+1} = 1$  and others are 0, which implies that the columns of  $B$  should be orthogonal to  $p_{k+1}$ . In this case, we take

$$B = [p_1, p_2, \dots, p_k], \quad x = p_{k+1}$$

The sup inf case is similar to the above. □

**Theorem 1.8** (Sturm separation theorem [3]). Consider a symmetric matrix  $A = A^T$ , we let

$$\{A_r = (a_{ij}) : i, j = 1, \dots, r; r = 1, 2, \dots, m\}$$

be the set of subset square matrices of  $A$ . Then:

$$\lambda_{k+1}(A_{i+1}) \leq \lambda_k(A_i) \leq \lambda_k(A_{i+1}) \quad (1.4)$$

**Theorem 1.9** (Sums of quadratic forms [3]). Suppose  $A$  is a symmetric matrix and  $x_1, \dots, x_k$  are  $k$  orthonormal vectors, that is,  $x_i^T x_j = \delta_{ij}$  where  $\delta_{ij}$  is the Kronecker delta. Then,

$$\sup_{x_1, \dots, x_k} \sum_{i=1}^k x_i^T A x_i = \sum_{i=1}^k \lambda_i$$

where  $\lambda_i$  is the  $i^{th}$  largest eigenvalue of  $A$ . The equality is attained when  $x_i$ 's are the corresponding eigenvectors to  $\lambda_i$ 's.

*Proof.* Let  $T \Lambda T^T$  be spectral decomposition of  $A$  and  $y_i = T^T x_i$ . We have  $y_i^T y_j = \delta_{ij}$ . WLOG, suppose  $y_i$  is  $p$ -dimension. Then LHS without sup can be written as

$$\sum_{i=1}^k \sum_{j=1}^p \lambda_j y_{ij}^2 = \sum_{j=1}^p \left( \sum_{i=1}^k y_{ij}^2 \right) \lambda_j$$

But  $\left( \sum_{i=1}^k y_{ij}^2 \right) \leq 1 \quad \forall j$  and  $\sum_{j=1}^p \left( \sum_{i=1}^k y_{ij}^2 \right) = k$  give us the desired result. □

**Theorem 1.10** (Cauchy-Schwarz inequality [3]). Let  $A = B^T B$  be a Gram matrix, then

$$1. \quad (x^T A y)^2 \leq (x^T A x)(y^T A y)$$

$$2. (x^T y)^2 \leq (x^T A x)(y^T A^{-1} y)$$

Thus,

$$\sup_x \frac{(u^T x)^2}{x^T A x} = u^T A^{-1} u$$

where  $A$  positive definite and  $x = A^{-1} u$ .

*Proof.* The proof is immediately obtained by letting

$$\tilde{x} = A^{1/2} x, \tilde{y} = A^{1/2} y$$

and

$$\tilde{x} = A^{1/2} x, \tilde{y} = A^{-1/2} y$$

□

## 1.5 Projection operator

Consider the classical setting  $y = X\beta + \epsilon$  where  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\beta \in \mathbb{R}^p$  and  $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$ . The projection operator  $P_X$  and  $Q_X$  are defined as

$$P_X = X(X^T X)^{-1} X^T, Q_X = I - P_X$$

In the general case when  $X$  is not invertible, then

$$P_X = X(X^T X)^- X^T$$

where  $(X^T X)^-$  is the g-inverse. Moreover, if  $\mathcal{V}$  is a vector space, then  $P_{\mathcal{V}}$  represents the projection operator onto  $\mathcal{V}$ . Unless otherwise specified, I use  $P$  instead of  $P_X, Q_X$  for simplicity. An interesting definition of projection operator in **finite dimension spaces** equipped with special inner product is given below.

**Definition 1.1** (Projection operator, [3]). Suppose  $\mathbb{R}^n$  is equipped with the inner product

$$\langle x, y \rangle = y^T \Sigma x$$

where  $x, y \in \mathbb{R}^n$  and  $\Sigma$  is a positive definite matrix. Let  $\mathcal{V}$  be a subspace of  $\mathbb{R}^n$  spanned by a square matrix  $P$ , then we say  $P$  is the projection operator on  $\mathcal{V}$  if the following holds

1.  $P$  is idempotent:

$$P^2 = P$$

2.  $\Sigma P$  is symmetric:

$$(\Sigma P)^T = \Sigma P$$

In the special case  $\Sigma + I$ , it is equivalent to say  $P$  is symmetric.

**Theorem 1.11** (Adding regressors). Let  $X = [X_1 \ X_2]$  be a  $n \times p$  matrix and  $X_1$  is  $n \times p_1$  while  $X_2$  is  $n \times p_2$  such that  $p = p_1 + p_2$ . The projection operator onto  $X$  can be written as

$$P_X = P_{X_1} + Q_{X_1} X_2 (X_2^T Q_{X_1} X_2)^{-1} X_2^T Q_{X_1}$$

*Proof.* For any vector  $\mathbf{m} \in \mathbb{R}^n$ , we have

$$\begin{aligned} P_X \mathbf{m} &= P_X P_{X_1} \mathbf{m} + P_X Q_{X_1} \mathbf{m} \\ &= P_{X_1} \mathbf{m} + P_X Q_{X_1} \mathbf{m} \\ &=_{(*)} P_{X_1} \mathbf{m} + P_{\mathcal{C}(Q_{X_1}) \cap \mathcal{C}(X)} \mathbf{m} \\ &= P_{X_1} \mathbf{m} + P_{\mathcal{C}(Q_{X_1}) \cap [\mathcal{C}(X_1) \cup \mathcal{C}(X_2)]} \mathbf{m} \\ &= P_{X_1} \mathbf{m} + P_{\mathcal{C}(Q_{X_1}) \cap \mathcal{C}(X_2)} \mathbf{m} \end{aligned}$$

(\*) is because  $\mathcal{C}(X) = \mathcal{C}(P_X) = \mathcal{C}(P_{X_1}) + \mathcal{C}(Q_{X_1}) \cap \mathcal{C}(X_2)$ . And the result follows from

$$P_{\mathcal{C}(Q_{X_1}) \cap \mathcal{C}(X_2)} = Q_{X_1} X_2 (X_2^T Q_{X_1} X_2)^{-1} X_2^T Q_{X_1}$$

or equivalently,

$$\mathcal{C}(Q_{X_1}) \cap \mathcal{C}(X_2) = \mathcal{C}(P_{(Q_{X_1} X_2)})$$

□

**Example 1.2** (Gram-Schmidt process). Let  $v, u$  be two  $p$ -dimensional vectors, then the projection operator associated with  $u$  is

$$P_u = \frac{uu^T}{\|u\|_2^2}$$

The projection of  $v$  onto  $u$  is

$$\text{proj}_u(v) = P_u(v) = \frac{\langle u, v \rangle}{\|u\|_2^2} u$$

Suppose we have  $k$  linearly independent vectors  $v_1, \dots, v_k$  and we want to orthogonalize them so

that they form an orthonormal basis of  $\mathbb{R}^k$ . Then the **Gram-Schmidt process** works as follows

$$\begin{aligned}
 u_1 &= v_1 & e_1 &= \frac{u_1}{\|u_1\|_2} \\
 u_2 &= v_2 - \text{proj}_{u_1}(v_2) & e_2 &= \frac{u_2}{\|u_2\|_2} \\
 u_3 &= v_3 - \text{proj}_{u_1}(v_3) - \text{proj}_{u_2}(v_3) & e_3 &= \frac{u_3}{\|u_3\|_2} \\
 &\dots & &\dots \\
 u_k &= v_k - \sum_{j=1}^{k-1} \text{proj}_{u_j}(v_k) & e_k &= \frac{u_k}{\|u_k\|_2}
 \end{aligned}$$

Finally,  $\{u_1, \dots, u_k\}$  forms an orthonormal basis of  $\mathbb{R}^k$ .

## 1.6 Exercises

1. (This formula comes from Hua Zhou's Biostat 257 HW1 for 2020 Spring, question 5) Let  $U$  and  $V$  be 2 matrices such that  $UV^T$  has the same dimension as  $A$ , show that

$$|A + UV^T| = |A| |I + V^T A^{-1} U|$$

This formula is useful for evaluating the density of a multivariate normal with covariance matrix  $A + UV^T$ .

2. (Binomial inversion theorem) Show that if  $A, B, C$  and  $D$  are conformable matrices, and all indicated inverses exist,

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}$$

Hence or otherwise, show Sherman-Morrison formula

$$(A + ab^T)^{-1} = A^{-1} - \frac{A^{-1}ab^T A^{-1}}{1 + b^T A^{-1}a}$$

and Woodbury identity

$$(A + UV)^{-1} = A^{-1} - A^{-1}U(I + V^T A^{-1}U)^{-1}V^T A^{-1}$$

This is useful when solving the Henderson's equation in linear mixed models.

3. Show the following:

(a) If  $A \in \mathbb{R}^{n \times m}$ , then  $\text{rank}(A) \leq \min(m, n)$ .

(b)  $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$ .

- (c)  $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$ .
- (d) If  $A, B$  are invertible, then  $\text{rank}(C) = \text{rank}(AC) = \text{rank}(CB) = \text{rank}(ACB)$ .
- (e) If  $\text{rank}(B) = \text{rank}(C)$ , then  $\mathcal{C}(ACB) = \mathcal{C}(AC)$ .
- (f) If  $A = ABA$ , then  $\text{rank}(A) = \text{rank}(AB) = \text{rank}(BA) \leq \text{rank}(B)$ .

4. (Key lemma) Show that

$$\text{rank}(A^T) = \text{rank}(A) = \text{rank}(AA^T) = \text{rank}(A^T A) \quad (1.5)$$

5. (Linear basis) If  $\text{rank}(A) = r$ , then show that there exists  $B$  and  $C$  such that

$$A = BC$$

where  $B$  has full column rank and  $C$  has full row rank.

6. (Eigenvalues and eigenvectors of symmetric matrices) The eigenvalues and eigenvectors of a matrix  $A$  can be obtained solving

$$Ax = (a + bi)x$$

where  $x \neq 0$  and  $a, b$  are real numbers and  $i$  is an indeterminate satisfying  $i^2 = -1$ . Show that  $b = 0$  if  $A$  is symmetric. **Hint:** Let  $Z = (A - (a + bi)I)$  and define

$$B = C^T \bar{C}$$

Then use the fact  $|B| = 0$ . Next, show  $x$ 's corresponding to distinct eigenvalues are orthogonal to each other.

- 7. Show that any matrix  $A$  can be written as a sum of a symmetric matrix and a skew matrix, and these summands are unique.
- 8. Prove that if  $x^T Ax = 0$  for all  $x$  and  $A = A^T$ , then  $A = 0$ .
- 9. Use the fact  $\mathcal{N}(A^T) = \mathcal{C}(A)^\perp$  to show if

$$PA^T A = QA^T A$$

then  $PA^T = QA^T$  for any conformable matrices  $P$  and  $Q$ .

10. (Simultaneous diagonalization) Let  $A$  and  $B$  be two  $n \times n$  positive definite matrices. If  $A - B \succ 0$ , show

- (a)  $\text{Tr}(A) > \text{Tr}(B)$ .
- (b)  $B^{-1} \succ A^{-1}$ , this is the theorem we have proved.

(c)  $\det |A| > \det |B|$ .

11. Show that the sets of non-zero eigenvalues for  $AB$  and  $BA$  are the same for any two conformable matrices.
12. Let  $A$  be a symmetric matrix and  $B = CA^{-1}C$ , describe the relationship between eigenvalues and eigenvectors of  $B$  and  $A$ .
13. Let  $\mathbf{1}_n$  be the  $n \times 1$  vector with all entries equal to 1 and let  $B$  be a  $n \times n$  matrix given by

$$B = \frac{2b}{2b-1}I_n - \frac{\mathbf{1}_n\mathbf{1}_n^T}{(2n-1)}$$

and  $b > 1/2$ .

- (a) Is  $B$  always positive definite?
- (b) Derive  $\text{Tr}(B)$ .
- (c) Show that the maximum eigenvalue of  $B$  is  $\frac{2b}{2b-1}$ .
14. (Moore-Penrose pseudo-inverse) Show that the Moore-Penrose pseudo-inverse  $A^+$  exists and is unique.
15. (More properties of  $g$ -inverse) Verify:
  - (a)  $\text{rank}(A) = \text{rank}(AA^-) = \text{rank}(A^-A)$ .
  - (b)  $\text{rank}(A) \leq \text{rank}(A^-)$ .
  - (c)  $\text{rank}(A) = \text{rank}(AA^+) = \text{rank}(A^+A)$ .
  - (d) If  $A$  has full column rank, then  $A^+ = (A^T A)^{-1}A$ .
  - (e) If  $A$  has full row rank, then  $A^+ = A^T(AA^T)^{-1}$ .
16. (Householder transformation) Let  $v$  be a nonzero vector. The **Householder transformation matrix** is defined by

$$H_v = I - \frac{2vv^T}{v^T v}$$

- (a) Find the determinant and all eigenvalues of such a matrix.
- (b) Show that if  $x \neq 0$ , then there is a Householder matrix such that

$$Hx = \|x\|_2 e_1$$

where  $e_1^T = (1, 0, 0, \dots, 0)$ .

17. (Norm-preserving mapping) Show that  $\|x\|_2 = \|y\|_2$  if and only if there is an orthogonal matrix  $T$  such that  $Tx = y$ .

18. (Canonical correlation) Suppose  $y$  is a univariate random vector with variance  $a^2$ , and  $X$  is a  $p \times 1$  random vector with covariance matrix  $V$  and

$$\text{Cov}(y, X) = W$$

where  $W$  is a  $p$ -dimensional vector. Use the generalized Cauchy-Schwarz inequality [1.10](#) to answer the following.

- (a) If  $b$  is any nonzero vector, what is the maximum correlation of  $b^T X$  with  $y$ ?
- (b) What choice of  $b$  will ensure that the maximum is attained?



## 2 Distribution Theory

This section is mostly based on the first half part of Biostat 250A where we discussed in depth the non-central  $\chi^2$  and  $t$  distribution with applications to linear models. Moment generating functions and characteristic functions are also included.

### 2.1 Multivariate normal distribution

In this subsection, we discuss basic properties of multivariate normal distribution.

**Definition 2.1** (Multivariate normal). A  $p$ -dimensional random vector  $Y$  is said to have a **multivariate normal distribution** with parameters  $(\mu, \Sigma)$  iff its moment generating function is

$$\Psi_Y(t) = \mathbb{E}e^{t^T y} = \exp\left(\mu^T t + \frac{1}{2}t^T \Sigma t\right) \quad (2.1)$$

We write

$$Y \sim \mathcal{N}_p(\mu, \Sigma)$$

and its density, if exists, is

$$f_Y(y) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left(-\frac{1}{2}(y - \mu)^T \Sigma^{-1}(y - \mu)\right)$$

**Lemma 2** (Alternative definition). *By Cramer-Wold device, we have*

$$Y \sim \mathcal{N}_p(\mu, \Sigma) \text{ iff } a^T y \sim \mathcal{N}_1(a^T \mu, a^T \Sigma a) \quad \forall a$$

**Lemma 3** (Linear transformation). *Let  $Y \sim \mathcal{N}(\mu, \Sigma)$ , then*

$$AY \perp\!\!\!\perp BY$$

*if and only if  $A\Sigma B = 0$  where  $A$  and  $B$  are two matrices of proper dimension.*

**Example 2.1** (Conditional multivariate normal distribution). Suppose

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} \sim \mathcal{N}_p\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}\right)$$

Set

$$\begin{pmatrix} Y_3 \\ Y_4 \end{pmatrix} = \begin{pmatrix} I & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & I \end{pmatrix} \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$$

Then we have

$$\text{Cov}(Y_3, Y_4) = \begin{pmatrix} \Sigma_{11 \cdot 2} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$$

where  $\Sigma_{11:2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ . Thus,  $Y_3 \perp\!\!\!\perp Y_2$  and we have

$$Y_1|Y_2 = y_2 \sim \mathcal{N}(\mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(y_2 - \mu_2), \Sigma_{11:2}) \quad (2.2)$$

## 2.2 Non-central distribution

In this subsection, we discuss non-central  $t$ ,  $\chi^2$  and  $F$  distribution and their statistical properties. Some results are from [3] and the derivation of the density of non-central  $\chi^2$  distribution is from Dr. Wong's 250B notes.

**Definition 2.2** (Non-central student  $t$ ). Let  $Y \sim \mathcal{N}(\mu, \sigma^2)$  be a univariate normal variable and  $\frac{X}{\sigma^2} \sim \chi_k^2(0)$  where  $\chi_k^2(0)$  is the usual central  $\chi^2$  distribution. Then we say

$$T = \frac{Y}{\sqrt{X/k}}$$

follows a **non-central student  $t$**  distribution with  $k$  degrees of freedom and non-centrality parameter  $\delta = \mu/\sigma$  and write  $T \sim t_k(\delta)$ . In the special case when  $Y \sim \mathcal{N}(\mu, 1)$ , we have  $\delta = \mu$  and write  $T \sim t_k(\mu)$ . The density of  $T$  is

$$f_t(t) = \frac{k^{\frac{k}{2}}}{\Gamma(\frac{k}{2})} \frac{\exp(-\frac{\delta^2}{2})}{(k+t^2)^{\frac{k+1}{2}}} \sum_{s=0}^{\infty} \left( \Gamma(\frac{k+s+1}{2}) \left(\frac{\delta^s}{s!}\right) \left(\frac{2t^2}{k+t^2}\right)^{\frac{s}{2}} \right)$$

Note there is a multivariate version of  $t$  distribution, see definition 7.3.

**Definition 2.3** (Non-central  $\chi^2$ ). Let  $X_i \sim \mathcal{N}(\mu_i, 1), i = 1, 2, \dots, k$  and  $X_i \perp\!\!\!\perp X_j$  for  $i \neq j$ . Define

$$\delta^2 = \sum_{i=1}^n \mu_i^2 = \|\mu\|_2^2$$

Then we say

$$Y = \sum_{i=1}^n X_i^2$$

has a non-central  $\chi^2$  distribution with  $n$  degrees of freedom and non-centrality parameter  $\delta^2$  and write

$$Y \sim \chi_n^2(\delta^2)$$

**Example 2.2** (Density of non-central  $\chi^2$ ). Suppose  $X \sim \mathcal{N}_n(\mu, I)$  where  $\mu^T = (\mu_1, \dots, \mu_n)$ . Let  $A$  be an orthogonal matrix derived from Gram-Schmidt process with the first row equals to  $\frac{\mu^T}{\|\mu\|_2}$ . Define

$$W = AX, Y = X^T X$$

Then

$$W \sim \mathcal{N}_n \left( \begin{pmatrix} \|\mu\| \\ 0 \\ \dots \\ 0 \end{pmatrix}, I \right)$$

and

$$\begin{aligned} Y &= X^T X \\ &= X^T A^T A X \\ &= W^T W \\ &= \underbrace{w_1^2}_{\chi_1^2(\|\mu\|_2^2)} + \underbrace{\sum_{i=2}^n w_i^2}_{\chi_{n-1}^2(0)} \end{aligned}$$

Thus, we can re-write a  $\chi_n^2(\delta^2)$  variable as

$$\begin{aligned} \chi_n^2(\delta^2) &= \chi_1^2(\delta^2) + \chi_{n-1}^2(0) \\ &= V + U \end{aligned} \tag{2.3}$$

and also  $V \perp U$ . Let  $Z \sim \mathcal{N}(\delta^2, 1)$ , then by definition,

$$\begin{aligned} F_V(v) &= \mathbb{P}(V \leq v) = \mathbb{P}(-\sqrt{v} \leq Z \leq \sqrt{v}) \\ &= \int_{-\sqrt{v}}^{\sqrt{v}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z-\delta^2)^2} dz \end{aligned}$$

Thus, the density of  $V$  is

$$\begin{aligned} f_V(v) &= \frac{d}{dv} F_V(v) \\ &= \frac{1}{\sqrt{2\pi}} \left( \frac{e^{-\frac{1}{2}(\sqrt{v}-\delta^2)^2}}{2\sqrt{v}} + \frac{e^{-\frac{1}{2}(-\sqrt{v}-\delta^2)^2}}{2\sqrt{v}} \right) \\ &= \frac{1}{2\sqrt{2\pi v}} e^{-\frac{1}{2}(v+\delta^4)} \left( e^{\delta^2\sqrt{v}} + e^{-\delta^2\sqrt{v}} \right) \\ &= \frac{1}{2\sqrt{2\pi v}} e^{-\frac{1}{2}(v+\delta^4)} \left\{ \sum_{i=0}^{\infty} \frac{(\delta^2\sqrt{v})^{2i}}{(2i)!} \right\} \end{aligned} \tag{2.4}$$

and  $\frac{(\delta^2\sqrt{v})^{2i}}{(2i)!}$ 's are derived from Taylor's expansion of  $e^x$ . The joint density of  $V$  and  $U$  is

$$f_{V,U}(v, u) = f_V(v) f_U(u)$$

with

$$f_U(u) = \frac{e^{-\frac{u}{2}} u^{\frac{n-3}{2}}}{2^{\frac{n-1}{2}} \Gamma(\frac{n-1}{2})}$$

Next, let  $Z_1 = U + V$  and  $Z_2 = \frac{U}{U+V}$  or equivalently,  $U = Z_1 Z_2$  and  $V = Z_1(1 - Z_2)$ . The Jacobian for this transformation is given by

$$J = \left\| \begin{pmatrix} Z_2 & Z_1 \\ 1 - Z_2 & -Z_1 \end{pmatrix} \right\| = Z_1$$

So,

$$\begin{aligned} f_{Z_1, Z_2}(z_1, z_2) &= f_V(z_1(1 - z_2)) f_U(z_1 z_2) \times z_1 \\ &= \text{const.} \times \exp\left(-\frac{1}{2}(z_1 + \delta^4)\right) z_1^{\frac{n-2}{2}} (1 - z_2)^{-\frac{1}{2}} \left\{ \sum_{i=0}^{\infty} (1 - z_2)^i \frac{\delta^{4i} z^i}{(2i)!} \right\} \end{aligned}$$

Since  $\int_0^1 (1 - x)^{i-\frac{1}{2}} dx = (i + \frac{1}{2})!$ , thus,

$$\begin{aligned} f_{Z_1}(z) &= \int_0^1 f_{Z_1, Z_2}(z, z_2) dz_2 \\ &= \text{const.} \times e^{-\frac{1}{2}(z+\delta^4)} z^{\frac{n-2}{2}} \left\{ \sum_{i=0}^{\infty} \frac{\delta^{4i} z^i}{(2i)!(i + \frac{1}{2})!} \right\} \end{aligned}$$

Recall  $\int_0^{\infty} x^{a-1} e^{-bx} dx = \Gamma(a)/b^a$ , so we get

$$f_{Z_1}(z) = \sum_{i=0}^{\infty} \underbrace{\frac{e^{-\lambda/2} (\frac{\lambda}{2})^i}{i!}}_{q_i} \underbrace{\left( \frac{1}{\Gamma(\frac{n+2i}{2})} \left(\frac{1}{2}\right)^{\frac{n+2i}{2}} z^{\frac{n+2i}{2}-1} e^{-\frac{z}{2}} \right)}_{\text{density of a } \chi_{n+2i}^2(0)} \quad (2.5)$$

In conclusion, we can re-write the density of  $Y \sim \chi_n^2(\delta^2)$  as

$$f_Y(y) = \sum_{i=0}^{\infty} q_i f_{\chi_{n+2i}^2(0)}(y) \quad (2.6)$$

and  $q_i$ 's are called **Poisson weights**.

**Definition 2.4** (Non-central  $F$ ). Let  $X \sim \chi_n^2(\delta^2)$  and  $Y \sim \chi_m^2(0)$  and assume  $X \perp\!\!\!\perp Y$ , then

$$F = \frac{X/n}{Y/m}$$

is said to have a non-central  $F$  distribution with parameter  $(\delta^2, n, m)$  and we write

$$F \sim F_{n,m}(\delta^2)$$

## 2.3 Basic results for quadratic forms

Suppose  $y \sim \mathcal{N}_p(0, I)$  and let

$$Q = y^T A y$$

where  $A^T = A$ . What is the distribution of  $Q$ ?

**Theorem 2.1** (Fundamental theorem). WLOG, let  $A$  be symmetric.

$$Q \sim \chi_r^2(0) \text{ iff } A^2 = A, \text{rank}(A) = r$$

*Proof.* It can be shown (exercise) the moment generating function (mgf) of  $Q$  is

$$\begin{aligned} \Psi_Q(t) &= \mathbb{E} e^{tQ} \\ &= \det |I - 2tA|^{-\frac{1}{2}} \\ &= \det |I - 2tD|^{-\frac{1}{2}} \\ &= \prod_{i=1}^p (1 - 2t\lambda_i)^{-\frac{1}{2}} \end{aligned} \tag{2.7}$$

where  $A = TDT^T$  is the spectral decomposition of  $A$ . If  $A^2 = A$ , then  $\lambda_i = 0$  or  $1$ , which implies

$$\Psi_Q(t) = (1 - 2t)^{-r/2}$$

But this is the mgf of  $\chi_r^2(0)$ . On the other hand, if the above holds, then

$$\prod_{i=1}^p (1 - 2\lambda_i t)^{-1/2} = (1 - 2t)^{-r/2} \quad \forall t$$

This means  $\lambda_i = 0$  for  $(p - r)$  different  $i$ 's and  $\lambda_i = 1$  for  $r$  different  $i$ 's.

□

**Example 2.3** (Sample variance). Suppose  $y \sim \mathcal{N}_p(0, \sigma^2 I)$ . Let

$$Q = y^T \left( I - \frac{\mathbf{1}\mathbf{1}^T}{n} \right) y$$

Then by fundamental theorem,

$$Q/\sigma^2 \sim \chi_r^2$$

where  $r = \text{Tr}(I - \frac{\mathbf{1}\mathbf{1}^T}{n}) = n - 1$ .

**Example 2.4** (General multivariate normal). Suppose  $y \sim \mathcal{N}_p(0, \Sigma)$ , let

$$W = \Sigma^{-\frac{1}{2}} y \sim \mathcal{N}_p(0, I)$$

and

$$Q = y^T A y = W^T \Sigma^{1/2} A \Sigma^{1/2} W$$

Then by fundamental theorem,

$$Q \sim \chi_r^2 \text{ iff } A \Sigma A = A, \text{rank}(A) = r$$

**Lemma 4** (Craig's independence lemma). *Let  $y \sim \mathcal{N}(0, I)$  and  $A^2 = A, B^2 = B$ . Then*

$$y^T A y \perp\!\!\!\perp y^T B y \text{ iff } AB = 0$$

*Proof.* If  $AB = 0$ , then  $Ay \perp\!\!\!\perp By$ . Since  $A, B$  are idempotent, we have

$$y^T A y = y^T A^2 y \perp\!\!\!\perp y^T B^2 y = y^T B y$$

If  $y^T A y \perp\!\!\!\perp y^T B y$ , then

$$y^T (A + B) y \sim \chi^2$$

By fundamental theorem,

$$(A + B)^2 = A + B$$

Thus,  $2AB = 0$ . □

**Theorem 2.2** (Hogg-Craig). Let  $y \sim \mathcal{N}(\mu, I)$  and  $Q_i = y^T P_i y$  for  $i = 1, 2$ . If  $Q_i \sim \chi_{r_i}^2(0)$  and  $Q_1 - Q_2 \geq 0, r_1 > r_2$ , then

$$Q_1 - Q_2 \perp\!\!\!\perp Q_2, Q_1 - Q_2 \sim \chi_{r_1 - r_2}^2(0)$$

*Proof.*  $0 \leq Q_1 - Q_2 = y^T (P_1 - P_2) y \forall y$ . In particular, if  $y \in \mathcal{N}(P_1)$ ,

$$0 \leq y^T (-P_2) y \leq 0$$

or

$$y^T P_2^T P_2^T y = 0$$

Thus,  $\mathcal{N}(P_1) \subset \mathcal{N}(P_2)$ . For any  $y, y^T P_2 (I - P_1) y = 0$  since  $(I - P_1) \in \mathcal{N}(P_1)$ . Therefore,

$$(P_1 - P_2)^2 = P_1^2 - P_1 P_2 - P_2 P_1 + P_2^2 = P_1 - P_2$$

and  $\text{rank}(P_1 - P_2) = \text{Tr}(P_1 - P_2) = r_1 - r_2$ . The result follows from fundamental theorem and independence lemma. □

## 2.4 General results for quadratic forms

**Lemma 5.** If  $y \sim \mathcal{N}(0, \Sigma)$ , then the mgf of  $Q = y^T A y$  is

$$\frac{1}{\det |I - 2t\Sigma A|^{1/2}} \text{ or } \frac{1}{\det |I - 2t\Sigma^{-1/2} A \Sigma^{1/2}|^{1/2}} \text{ or } \frac{1}{\det |I - 2tA\Sigma|^{1/2}}$$

*Proof.* Use the fact  $|I - AB| = |I - BA|$ . □

**Theorem 2.3** (Independence of quadratic forms). Let  $y \sim \mathcal{N}(0, \Sigma)$ ,  $Q_i = y^T A_i y$ ,  $i = 1, 2$ . Then

$$Q_1 \perp\!\!\!\perp Q_2 \text{ iff } A_1 \Sigma A_2 = 0$$

*Proof.* If  $A_1 \Sigma A_2 = 0$ , then

$$\begin{aligned} \Psi_{Q_1, Q_2}(t_1, t_2) &= \mathbb{E} e^{t_1 Q_1 + t_2 Q_2} \\ &= \int e^{t_1 Q_1 + t_2 Q_2} \frac{1}{2\pi |\Sigma|^{1/2}} e^{-\frac{1}{2} y^T \Sigma^{-1} y} dy \\ &= \frac{1}{2\pi |\Sigma|^{1/2}} \int e^{-\frac{1}{2} y^T [\Sigma^{-1} - 2t_1 A_1 - 2t_2 A_2] y} dy \\ &= \frac{|\Sigma^{-1} - 2t_1 A_1 - 2t_2 A_2|^{-1/2}}{|\Sigma|^{1/2}} \\ &= \frac{1}{|I - 2t_1 A_1 \Sigma - 2t_2 A_2 \Sigma|^{1/2}} \\ &\stackrel{(*)}{=} \frac{1}{|I - 2t_1 A_1 \Sigma|^{1/2} |I - 2t_2 A_2 \Sigma|^{1/2}} \\ &= \Psi_{Q_1}(t_1) \Psi_{Q_2}(t_2) \end{aligned}$$

where  $(*)$  follows from the fact  $A_1 \Sigma A_2 = 0$ . The other direction follows similarly. □

Now let  $Y \sim \mathcal{N}_p(\mu, \Sigma)$ ,  $\Sigma \succ 0$  and we want to find the distribution of  $Q = y^T A y$ , given that  $\text{rank}(A) = p$ . Write

$$Q = y^T \Sigma^{-1/2} T T^T \Sigma^{1/2} A \Sigma^{1/2} T T^T \Sigma^{-1/2} y$$

where  $T T^T = T^T T = I$  and  $T^T \Sigma^{1/2} A \Sigma^{1/2} T = \Lambda$ , i.e., the spectral decomposition of  $\Sigma^{1/2} A \Sigma^{1/2}$ . Hence,

$$Q = w^T \Lambda w = \sum_{i=1}^p \lambda_i w_i^2$$

where  $\lambda_i$ 's are eigenvalues of  $\Sigma^{1/2} A \Sigma^{1/2}$  and

$$w \sim \mathcal{N}_p(T^T \Sigma^{-1/2} \mu, I)$$

Thus,  $Q = \sum_{i=1}^r \lambda_i w_i^2$  is a weighted sum of non-central  $\chi_1^2(\delta_i^2)$  variables where

$$\delta_i^2 = (t_i^T \Sigma^{-1/2} \mu)^2$$

and  $t_i$  is the  $i^{th}$  column of  $T$ ,  $r$  is the rank of  $A$ .

**Example 2.5** ( $\mu = 0, \Sigma = I, A^2 = A$ ). Then

$$Q =_d \sum_i \lambda_i \chi_1^2(0)$$

where  $\lambda_i$ 's are non-zero eigenvalues of  $A$ . In other words,

$$Q \sim \chi_r^2(0)$$

**Example 2.6** ( $\mu = 0, A = \Sigma^{-1}$ ). Then

$$\begin{aligned} Q &= y^T A y \\ &= y^T \Sigma^{-1/2} \Sigma^{1/2} A \Sigma^{1/2} \Sigma^{-1/2} y \\ &= \left( \Sigma^{-1/2} y \right)^T \left( \Sigma^{-1/2} y \right) \end{aligned}$$

Hence,

$$Q \sim \chi_p^2(0)$$

**Example 2.7** ( $\mu \neq 0, A = \Sigma^{-1}$ ). In this case,

$$\delta_i^2 = \mu^T \Sigma^{-1/2} t_i t_i^T \Sigma^{-1/2} \mu$$

and thus,

$$\sum_{i=1}^p \delta_i^2 = \mu^T \Sigma^{-1/2} \left( \sum_{i=1}^p t_i t_i^T \right) \Sigma^{-1/2} \mu$$

and  $\sum_{i=1}^p t_i t_i^T = T T^T = I$ . Since  $T^T \Sigma^{1/2} A \Sigma^{1/2} T = I$ ,  $\lambda_1 = \dots = \lambda_p = 1$ . Hence,

$$Q \sim \chi_p^2(\mu^T \Sigma^{-1} \mu)$$

**Example 2.8** ( $\mu \neq 0, \Sigma = I, A^2 = A$ ). In this case, we have

$$Q \sim \chi_r^2 \left( \sum_{i=1}^r \delta_i^2 \right)$$

where

$$\sum_{i=1}^r \delta_i^2 = \mu^T A \mu$$



## 2.5 Fisher-Cochran's theorem

Fisher-Cochran's theorem is a powerful tool to prove many fundamental results in least square theory. To show Fisher-Cochran's theorem, we need the following two lemmas: Loynes lemma and Marsaglia-Garaybill lemma.

**Lemma 6** (Loynes). *If  $M^2 = M = M^T$ ,  $P^T = P \succeq 0$  and  $I - M - P \succeq 0$ . Then*

$$PM = MP = 0$$

*Proof.* Let  $y \in \mathbb{R}^n$  and  $z = My$ , then  $z^T = y^T M^T M y$  or

$$z^T z = y^T M M M y = z^T M z$$

Thus,

$$z^T (I - M) z = 0$$

Next, by assumption,

$$0 \leq z^T (I - M - P) z = -z^T P z \leq 0$$

Thus,

$$z^T P z = 0 \Rightarrow P z = 0 \text{ or } P M y = 0 \forall y$$

which suggests

$$PM = MP = 0$$

□

The next lemma was firstly proved by Garabill and Marsaglia in 1957 [4], then in 1964, K. S. Banerjee provided a cleaner proof.

**Lemma 7** (Marsaglia-Garaybill, [4] and [5]). *Suppose for  $1 \leq i \leq k$ ,*

$$D_i = D_i^T$$

*Then **any two** of the following statements imply the third:*

1.  $D_i^2 = D_i, i = 1, 2, \dots, k.$
2.  $D_i D_j = 0 \forall i \neq j.$
3.  $D = \sum_{i=1}^k D_i$  is idempotent.

*Proof.* The proof has 3 parts.

- (1)(2)  $\Rightarrow$  (3): trivial.

- (3)(1)  $\Rightarrow$  (2): Since for any  $i$ ,  $D_i$  is symmetric, then  $D_i^2 = D_i$  implies

$$D_i \succeq 0 \quad \forall i$$

Thus,

$$D - D_i - D_j = \sum_{k \neq i, j} D_k \succeq 0$$

Next, by third condition,

$$D^2 = D, D^T = D \Rightarrow I - D \succeq 0$$

Hence

$$I - D_i - D_j = (I - D) + (D - D_i - D_j) \succeq 0$$

By Loynes' lemma where

$$M = D_i, P = D_j$$

we have

$$D_i D_j = D_j D_i = 0$$

for any  $i \neq j$ .

- (2)(3)  $\Rightarrow$  (1): By definition of eigenvalues,

$$D_i x = \lambda x \Rightarrow D D_i x = \lambda D x$$

or by second condition,

$$\lambda^2 x = D_i^2 x = \lambda D x$$

Thus, by third condition,  $\lambda = 0$  or  $1$ , which means  $D_i^2 = D_i$  is idempotent.

□

**Theorem 2.4** (Fisher-Cochran). Suppose

$$y \sim \mathcal{N}_p(0, I)$$

and

$$y^T y = \sum_{i=1}^k Q_i$$

where  $Q_i = y^T A_i y$ ,  $\text{rank}(A_i) = r_i$  and each  $A_i$  is, of course, symmetric. Then T.F.A.E.:

1.  $Q_i \perp Q_j \quad \forall i \neq j$ .
2.  $Q_i \sim \chi_{r_i}^2(0), i = 1, 2, \dots, k$ .
3.  $\sum_{i=1}^k r_i = p$ .

*Proof.* I will show that  $(i) \Rightarrow (ii) \Rightarrow (iii) \Rightarrow (i)$ .

- $(i) \Rightarrow (ii)$ :

$$Q_i \perp\!\!\!\perp Q_j \Rightarrow Q_1 \perp\!\!\!\perp \sum_{i=2}^k Q_i$$

or equivalently,

$$y^T A_1 y \perp\!\!\!\perp y^T \left( \sum_{i=2}^k A_i \right) y$$

Then by Craig's independence lemma 4 we have

$$A_1 \left( \sum_{i=2}^k A_i \right) = 0$$

Since  $\sum_{i=1}^k A_i = I$ ,

$$A_1(I - A_1) = 0 \Rightarrow A_1^2 = A_1 \Rightarrow A_i^2 = A_i, i = 1, \dots, k$$

Hence by fundamental theorem 2.1,

$$Q_i \sim \chi_{r_i}^2(0)$$

- $(ii) \Rightarrow (iii)$ : Note by fundamental theorem 2.1,

$$Q_i \sim \chi_{r_i}^2(0) \Rightarrow A_i^2 = A_i$$

Thus,

$$\sum_{i=1}^k r_i = \sum_{i=1}^k \text{Tr}(A_i) = \text{Tr}\left(\sum_{i=1}^k A_i\right) = \text{Tr}(I) = p$$

- $(iii) \Rightarrow (i)$ : This is the part that we need **Marsaglia-Garaybill's lemma** (or Loynes' lemma).

Let  $A = \sum_{i=2}^k A_i$ , then

$$A_1 + A = I$$

Let  $T$  be such that  $T^T = TT^T = I$  and

$$T^T A_1 T = \Lambda$$

where  $T \Lambda T^T$  is the spectral decomposition of  $A_1$ . Then

$$T^T A_1 T + T^T A T = T^T T = I$$

or

$$\Lambda + T^T AT = I$$

**KEY STEP:** but

$$\text{rank}(T^T AT) = \text{rank}(A) = \text{rank}\left(\sum_{i=2}^p A_i\right) \leq \sum_{i=2}^p \text{rank}(A_i) = p - r_1 \quad (2.8)$$

where the last equality follows from the third condition. Denote the non-zero diagonal entries of  $\Lambda$  as  $\lambda_1, \dots, \lambda_{r_1}$ , and the rest are 0's. Then the corresponding elements of  $T^T AT$  are  $1 - \lambda_1, \dots, 1 - \lambda_{r_1}$  and 1's respectively. However, inequality 2.8 means

$$1 - \lambda_i = 0 \quad \forall i = 1, 2, \dots, r_1$$

or equivalently,

$$\lambda_i = 1 \quad \forall i = 1, 2, \dots, r_1$$

Therefore,

$$A_1^2 = A_1 \Rightarrow A_i^2 = A_i, i = 1, 2, \dots, k$$

which means  $A_i \succeq 0$  for  $i = 1, \dots, k$ . Next,

$$I - A_i - A_j = \sum_{k \neq i, j} A_k \succeq 0$$

and hence by Loynes' lemma

$$A_i A_j = A_j A_i = 0 \quad \forall 1 \leq i \neq j \leq k$$

and these two steps can be derived directly by Marsaglia-Garaybill's lemma. But by Craig's independence lemma,

$$Q_i \perp\!\!\!\perp Q_j \text{ iff } A_i \Sigma A_j = 0$$

Thus,

$$Q_i \perp\!\!\!\perp Q_j \quad \forall 1 \leq i \neq j \leq k$$

□

## 2.6 Applications of Fisher-Cochran's theorem

**Example 2.9.** Suppose  $y \sim \mathcal{N}(0, I)$ . We have

$$y^T y = y^T \frac{\mathbf{1}\mathbf{1}^T}{n} y + y^T \left(I - \frac{\mathbf{1}\mathbf{1}^T}{n}\right) y$$

By Fisher-Cochran, we have

$$y^T \left( I - \frac{\mathbf{1}\mathbf{1}^T}{n} \right) \sim \chi_r^2(0), \quad r = n - 1$$

$$y^T \frac{\mathbf{1}\mathbf{1}^T}{n} y \sim \chi_1^2(0)$$

**Example 2.10** (Multiple regression). In regression problem:

$$y = X\beta + \epsilon$$

where  $\beta = (\beta_0, \beta_1, \dots, \beta_{p-1})$ . We want to test (**Omnibus test**)

$$H_0 : \beta_1 = \dots = \beta_{p-1} = 0$$

and note that  $\hat{\beta} = (X^T X)^{-1} X^T y$ ,  $\text{rank}(X) = p$ . Then

$$y^T y = \underbrace{y^T (I - P) y}_{\text{SSE}} + \underbrace{y^T \left( P - \frac{\mathbf{1}\mathbf{1}^T}{n} \right) y}_{\text{SSReg}} + y^T \frac{\mathbf{1}\mathbf{1}^T}{n} y \quad (2.9)$$

Then by Fisher-Cochran, the test statistic is

$$\frac{\text{SSReg}}{\text{SSE}} \frac{n-p}{p-1} \sim_{H_0} F_{p-1, n-p}$$

**Example 2.11** (One-way ANOVA). Suppose for  $i = 1, \dots, k$  and  $j = 1, \dots, n$ ,

$$y_{ij} = \mu + \tau_i + \epsilon_{ij}$$

$$\epsilon_{ij} \sim_{iid} \mathcal{N}(0, \sigma^2)$$

How do we test  $\tau_1 = \dots = \tau_k$ ? Note that

$$\sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y})^2 = \underbrace{\sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{i\cdot})^2}_{\text{SSE}} + \underbrace{n \sum_{i=1}^k (\bar{y}_{i\cdot} - \bar{y})^2}_{\text{SSTr}} \quad (2.10)$$

By Fisher-Cochran, we have

$$\frac{\text{SSE}}{\sigma^2} = \sum_{i=1}^k \sum_{j=1}^n (\epsilon_{ij} - \bar{\epsilon}_{i\cdot})^2 / \sigma^2 \sim_{H_0} \sum_{i=1}^k \chi_{n-1}^2$$

and  $\mathbb{E}(\frac{\text{SSE}}{(n-1)k}) = \sigma^2$ ,

$$\frac{\text{SSTr}}{\sigma^2} = \sum_{i=1}^k \sum_{j=1}^n (\bar{y}_{i\cdot} - \bar{y})^2 / \sigma^2 \sim_{H_0} \chi_{k-1}^2$$

and  $\mathbb{E}(\frac{\text{SSTr}}{k-1}) = \sigma^2 + n \frac{\sum_{i=1}^k (\tau_i - \bar{\tau})^2}{k-1}$ . For more on these, see **ANOVA mixed models**.

## 2.7 Exercises

1. (Reproducing property) Let  $X_i, i = 1, \dots, n$  be a sequence of normal variables with mean  $\mu_i$ . Let  $Y_1 = \sum_{i=1}^r X_i^2$  and  $Y_2 = \sum_{i=r+1}^n X_i^2$ . Show that for  $0 < t < 0.5$ ,

- (a) The mgf of  $Y_1$  is

$$\Psi_{Y_1}(t) = \frac{e^{\delta_1^2 t / (1-2t)}}{(1-2t)^{r/2}}$$

- (b) The mgf of  $Y = Y_1 + Y_2$  is

$$\Psi_Y(t) = \frac{e^{\lambda t / (1-2t)}}{(1-2t)^{n/2}}$$

where  $\lambda = \delta_1^2 + \delta_2^2$ ,  $\delta_1^2 = \sum_{i=1}^r \mu_i^2$  and  $\delta_2^2 = \sum_{i=r+1}^n \mu_i^2$ .

2. Prove formula 2.7.

3. (MVN when covariance is singular) Suppose  $Y$  has a  $p$ -dimensional multivariate normal distribution with mean  $m$  and covariance  $V$ . Let  $Y$  be partitioned into  $Y_1$  and  $Y_2$  so that  $Y^T = (Y_1^T, Y_2^T)$  and the dimension of  $Y_1$  and  $Y_2$  are, respectively,  $p_1$  and  $p_2$  with  $p_1 + p_2 = p$ . Find the distribution of  $Y_1$  given  $Y_2$  if  $V$  is **singular** and express the conditional distribution in terms of only  $m$ ,  $p_1$  and the four submatrices in the partitioned matrix  $V$  and an appropriate  $g$ -inverse.

4. (Non-central  $\chi^2$ ) Let  $m$  be a  $p$ -dimensional vector,  $y \sim \mathcal{N}_p(m, I)$  and  $A$  is idempotent of rank  $k$ . What is the distribution of  $(y - a)^T A (y - a)$ ?
5. (Density of non-central  $\chi^2$  and its mgf) Let  $p_m$  be the density of a central  $\chi^2$  distribution with  $m$  degrees of freedom, and for each non-negative  $s$ , let

$$q_j = \frac{(s/2)^j}{j!} \exp(-s/2)$$

for  $j = 0, 1, \dots$ . A random variable with density

$$h(z) = \sum_{j=0}^{\infty} q_j p_{m+2j}(z), \quad z > 0$$

is said to have a non-central  $\chi^2$  distribution with  $m$  degrees of freedom and non-centrality parameter  $s$ . Use the **power series expansion** of  $\exp(x)$  and **monotone convergence theorem** to find the moment generating function of such a random variable.

6. (Non-central  $F$ ) Find the expectation of a non-central  $F$  distribution with numerator and denominator degrees of freedom  $n$  and  $m$ , respectively, and non-centrality parameter  $\Phi$ .
7. Let  $X_1, \dots, X_{n_1}$  and  $Y_1, \dots, Y_{n_2}$  be ind random samples from  $\mathcal{N}(\mu_1, v_1^2)$  and  $\mathcal{N}(\mu_2, v_2^2)$  respectively. Let  $\bar{X}, \bar{Y}, S_1^2, S_2^2$  denote the respective sample means and variances, and let  $c$  be a fixed constant. Identify the distribution of each of the following statistics by finding a suitable value of  $k$ :

- (a)  $\frac{\sqrt{n_2}(\bar{Y}-c)}{kv_2}$ .
- (b)  $\frac{k\sqrt{n_1}(\bar{X}-c)}{S_1}$ .
- (c)  $\frac{k(X_1+X_2)}{|Y_1-Y_2|}$ .
- (d)  $\frac{k[(X_1-c)^2+(X_2-c)^2]}{S_2^2}$ .

8. Let  $(X_j, Y_j), j = 1, 2, \dots, n$  be a random sample from the bi-variate normal distribution with parameters  $m_1, m_2, v_1^2, v_2^2$  and correlation  $r$ .

- (a) If  $d$  is a fixed constant, find a constant  $k$  s.t.

$$T = \frac{k(\bar{X} - \bar{Y} - d)}{\sqrt{\sum_{i=1}^n (X_i - Y_i - \bar{X} + \bar{Y})^2}}$$

has a non-central  $t_m(s)$  distribution.

- (b) Express  $m$  and  $s$  as a function of the parameters and the constant  $d$ .
- (c) What is the expectation of  $T$ ?

9. (250A Midterm) (The ANOVA Theorem) Let  $V$  be a  $n \times n$  positive definite matrix, let  $y \sim \mathcal{N}(m, V)$ , let  $A$  be a  $n \times n$  symmetric matrix of rank  $r$  and let

$$A = A_1 + A_2 + \dots + A_k$$

where each  $A_i$  is symmetric of rank  $r_i, i = 1, 2, \dots, k$ . The following conditions may be defined:

- (i)  $A_i V = (A_i V)^2, i = 1, 2, \dots, k$ ;
- (ii)  $A_i V A_j = 0, 1 \leq i \neq j \leq k$ ;
- (iii)  $AV = (AV)^2$ .

Show that the following statements are all simultaneously true if and only if **any two** of (i), (ii) and (iii) are true:

- (a)  $y^T A_i y \sim \chi_{r_i}^2(m^T A_i m), i = 1, 2, \dots, k.$
- (b)  $y^T A_i y \perp\!\!\!\perp y^T A_j y, 1 \leq i \neq j \leq k.$
- (c)  $y^T A y \sim \chi_r^2(m^T A m).$

10. (250A Final) Suppose that  $Y \sim \mathcal{N}_p(0, I_p).$

- (a) Find the conditional distribution of  $Y$  given  $1^T Y = 0$  and show that the covariance matrix of the conditional distribution is always positive definite.
- (b) If  $a$  is a non-zero vector, find the conditional distribution of  $Y^T Y$  given  $a^T Y = 0$  and find its expectation.

11. (250A Final) If  $Y^T = (Y_1, Y_2) \sim \mathcal{N}_2(0, V)$  and the  $(i, j)^{th}$  element of  $V$  is  $v_{ij}$ , find the distribution of  $Y^T V^{-1} Y - Y_1^2 / v_{11}.$

12. (250A Final) Let  $x_1 = (1, 1, 1, 1, 1)^T, x_2 = (1, 1, 0, 0, 0)^T, \theta = (6, 6, 2, 2, 2)^T$  and suppose  $Y \sim \mathcal{N}_5(\theta, \sigma^2 I_5).$  Let  $V$  be the linear span of  $x_1$  and  $x_2$  and let  $\hat{Y}$  be the OP of  $Y$  onto  $\mathcal{V}$ . Find a constant  $K$  so that

$$\frac{K \|\hat{Y}\|_2^2}{\|Y - \hat{Y}\|_2^2}$$

has one of the distributions discussed in class. Make sure you completely specify the distribution and identify all the parameters in the distribution.



### 3 Theory of Least Square Estimation

In 250A, we talked about basics of linear models, properties of LS and G-LS estimators, constrained LSE, adding regressors, more on orthogonal projections, rank-deficient case and some hypothesis testing problems.

#### 3.1 Basics of linear models

Let the model be

$$\begin{aligned}y &= X\beta + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 I)\end{aligned}\tag{3.1}$$

where  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\beta \in \mathbb{R}^p$ ,  $\epsilon \in \mathbb{R}^n$ . This is the most classic linear model in statistics. The normal equation is defined as

$$X^T X \beta = X^T y$$

There is always a solution to the normal equation because

$$X^T y \in \mathcal{C}(X^T) = \mathcal{C}(X^T X)$$

Let

$$\hat{\beta} = (X^T X)^- X^T y$$

be a least square estimate of  $\beta$ . The  $g$ -inverse  $(X^T X)^-$  is not unique, however, the projection

$$\hat{y} = X \hat{\beta}$$

is indeed unique. Define  $\theta = X\beta$ , we have the famous **Gauss-Markov BLUE** theorem.

**Theorem 3.1** (BLUE). The least square estimate  $\hat{y}$  of  $\theta$  is BLUE. That is, if interest is in estimating  $c^T \theta$  where  $c$  is given, then  $c^T \hat{y}$  is the best unbiased linear estimator.

The proof (exercise) follows directly from the fact that  $\hat{y}$  is unbiased and can be re-written as the projection of  $y$  on  $X$ , i.e.,  $\hat{y} = P_X y$ .

#### 3.2 Properties of LS and G-LS estimators

**Example 3.1** (Distribution of  $\hat{\beta}, s^2$ ). Under assumption of normality, we have

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^T X)^{-1})\tag{3.2}$$

$$s^2 \sim \frac{1}{n-p} \chi_{n-p}^2(0)\tag{3.3}$$

where

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T y \\ s^2 &= \frac{1}{n-p} y^T Q_X y\end{aligned}$$

Further, look at  $Cov(Q_X y, \hat{\beta})$ , then we have  $Q_X y \perp \hat{\beta}$ . Hence,

$$s^2 \perp \hat{\beta} \quad (3.4)$$

**Example 3.2** (Weighted LS or G-LS). Suppose that

$$\epsilon \sim \mathcal{N}(0, \sigma^2 V)$$

where  $V \succ 0$  is a known covariance matrix. Then we transform the model to be:

$$V^{-1/2} y = V^{-1/2} X \beta + V^{-1/2} \epsilon$$

Hence, the G-LS of  $\beta$  is

$$\hat{\beta}_W = (X^T V^{-1} X)^{-1} X^T V^{-1} y \quad (3.5)$$

and the covariance is

$$Cov(\hat{\beta}) = (X^T V^{-1} X)^{-1} \quad (3.6)$$

### 3.3 Adding regressors

Let the design matrix be partitioned as

$$X = \begin{pmatrix} X_1 & X_2 \end{pmatrix}$$

If  $X_1 \perp X_2$ , then the LSE of  $\beta$  can be written as

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} (X_1^T X_1)^{-1} X_1^T y \\ (X_2^T X_2)^{-1} X_2^T y \end{pmatrix}$$

The question is, how to make them orthogonal? Let the  $G$ -model be:

$$\mathbb{E}y_G = X\beta + Z\gamma$$

and  $X \in \mathbb{R}^{n \times p}$ ,  $Z \in \mathbb{R}^{n \times t}$ , both are of full column rank.

**Lemma 8.** *If  $Z^T X = 0$ , then  $Z^T Q_X Z$  is non-singular.*

*Proof.* The proof is simple:  $Q_X Z = (I - P_X)Z = Z$  and  $Z$  is of full column rank. □

Next, re-write the model as

$$\begin{aligned} y &= X\beta + P_X Z\gamma + Q_X Z\gamma \\ &= X\alpha + Q_X Z\gamma \end{aligned}$$

where  $\alpha = \beta + (X^T X)^{-1} X^T Z\gamma$ . Then we have

$$\hat{\alpha} = (X^T X)^{-1} X^T y \quad (3.7)$$

$$= \hat{\beta}_G + (X^T X)^{-1} X^T Z\hat{\gamma}_G$$

$$\hat{\beta}_G = (X^T X)^{-1} X^T (y - Z\hat{\gamma}_G) \quad (3.8)$$

$$\hat{\gamma}_G = (Z^T Q_X Z)^{-1} Z^T Q_X y \quad (3.9)$$

The above argument also provides a different proof of theorem 1.11. On the other hand, we have apply theorem 1.11 to derive the estimators  $\hat{\alpha}$  and  $\hat{\gamma}_G$  directly.

**Example 3.3** (Residuals from the  $G$ -model). Let the prediction of  $y$  be

$$\hat{y}_G = X\hat{\beta}_G + Z\hat{\gamma}_G$$

Then the residual vector is

$$E_G = y - \hat{y}_G = Q_X (y - Z\hat{\gamma}_G)$$

The sum of square error is

$$\begin{aligned} \text{SSE}_G &= e_G^T e_G \\ &= y^T Q_X y - \hat{\gamma}_G^T Z^T Q_X y - \hat{\gamma}_G^T (Z^T Q_X y - Z^T Q_X Z\hat{\gamma}_G) \end{aligned}$$

The first term is  $s^2$ , the third term is indeed zero. The second term is no less than 0 since

$$\hat{\gamma}_G^T Z^T Q_X y = y^T Q_X Z (Z^T Q_X Z)^{-1} Z^T Q_X y \geq 0$$

In conclusion,  $\text{SSReg}_G$  cannot be smaller when we add regressors. However, we have (exercise) [Okay fine, actually I am too lazy to type all the intermediate steps]

$$\text{Cov}(\hat{\beta}_G) \geq \text{Cov}(\hat{\beta}) \quad (3.10)$$

### 3.4 More on projections

Suppose  $\mathcal{V}$  is a  $r$ -dimensional subspace of  $\Omega$ , which is  $n$ -dimension. Given  $y \in \Omega$ ,  $\hat{y}$  is an orthogonal projection of  $y$  onto  $\mathcal{V}$  if

- $\hat{y} \in \mathcal{V}$ .
- $y - \hat{y} \in \mathcal{V}^\perp$ .

Now the question is, how to find  $\hat{y}$ ? The answer is trivial: suppose  $\mathcal{V}$  has  $\{x_1, \dots, x_r\}$  as an orthonormal basis. Then we use Gram-Schmidt to extend it to  $n$ -dimension:

$$\{x_1, \dots, x_r, x_{r+1}, \dots, x_n\}$$

and this is an orthonormal basis for  $\Omega$ . Then:

$$\begin{aligned} y &= \sum_{i=1}^n \alpha_i x_i \\ &= \sum_{i=1}^r \alpha_i x_i + \sum_{j=r+1}^n \alpha_j x_j \\ &= \hat{y} + \sum_{j=r+1}^n \alpha_j x_j \\ y - \hat{y} &= \sum_{j=r+1}^n \alpha_j x_j \in \mathcal{V}^\perp \end{aligned}$$

**Lemma 9.** *The projection  $\hat{y}$  is unique.*

*Proof.* WLOG, suppose  $\hat{y}_1$  and  $\hat{y}_2$  are two distinct projections of  $y$  onto  $\mathcal{V}$ . Then,

$$y = \hat{y}_1 + \hat{y}_1^\perp = \hat{y}_2 + \hat{y}_2^\perp$$

Hence,

$$\hat{y}_1 - \hat{y}_2 = \hat{y}_2^\perp - \hat{y}_1^\perp$$

But LHS belongs to  $\mathcal{V}$  and RHS belongs to  $\mathcal{V}^\perp$  and  $\mathcal{V} \cap \mathcal{V}^\perp = \{0\}$ . Therefore,

$$\hat{y}_1 = \hat{y}_2$$

□

**Lemma 10.** *The  $\alpha_j$ 's are Fourier coefficients of  $y$  w.r.t.  $x_j$ 's, i.e.,*

$$\alpha_j = \langle x_j, y \rangle$$

*Proof.* See exercise.

□

Next, define

$$T = \begin{bmatrix} x_1 & x_2 & \cdots & x_r \end{bmatrix}$$

and

$$P = TT^T$$

Then  $P$  is an orthogonal projection (OP) matrix onto  $\mathcal{V}$ . This OP is unique. Another example concerns  $(X^T X)^-$ , which is not unique in general. Remarkably,

$$P = X(X^T X)^- X^T \quad (3.11)$$

is indeed UNIQUE (exercise).

**Lemma 11** (Two facts about OP matrix). *Let  $\mathcal{W}$  and  $\Omega$  be two vector space.*

1. *If  $\mathcal{W} \subset \Omega$ , then*

$$P_{\mathcal{W}}P_{\Omega} = P_{\Omega}P_{\mathcal{W}} = P_{\mathcal{W}} \quad (3.12)$$

2. *If  $\mathcal{W} \subset \Omega$ , then we have*

$$P_{\Omega} - P_{\mathcal{W}} = P_{\mathcal{W}^{\perp} \cap \Omega} \quad (3.13)$$

Part one is also known as **tower property of conditional expectation**. That is, if  $\mathcal{F}_1 \subset \mathcal{F}_2$  are two  $\sigma$ -algebras and  $X$  is a random variable from  $\mathcal{F}$  to  $\mathcal{B}(\mathbb{R})$ , then we have

$$\mathbb{E}(\mathbb{E}(X|\mathcal{F}_1)|\mathcal{F}_2) = \mathbb{E}(\mathbb{E}(X|\mathcal{F}_2)|\mathcal{F}_1)$$

But the above does not hold in general if  $\mathcal{F}_1$  and  $\mathcal{F}_2$  are not adapted. Part two follows from the fact that

$$\Omega = \mathcal{W} + \mathcal{W}^{\perp} \cap \Omega$$

provided  $\mathcal{W} \subset \Omega$ . The following theorem allows us to provide an analytical solution to the constrained least square estimation in next section.

**Theorem 3.2** (Key lemma for constrained LSE). Let  $\mathcal{W}$  and  $\Omega$  be two vector sub-space of  $\mathcal{V}$ . Suppose  $\mathcal{W} = \mathcal{N}(A) \cap \Omega$  for some  $A$ , then we have

$$\mathcal{W}^{\perp} \cap \Omega = \mathcal{C}(P_{\Omega}A^T) \quad (3.14)$$

As a by-product,

$$P_{\mathcal{W}^{\perp} \cap \Omega} = (P_{\Omega}A^T)(AP_{\Omega}A^T)^{-1}AP_{\Omega}$$

*Proof.* Here I will provide my own proof because Dr. Wong's proof is too abstract and I cannot fully understand some intermediate steps. My proof is similar to the proof of 1.11.

Firstly, note that

$$\mathcal{W}^{\perp} = \mathcal{N}(A)^{\perp} \bigoplus \Omega^{\perp}$$

The reason is that the operator  $(\cdot)^\perp$  can be regarded as taking complement of a set and by De Morgan's law we have  $(A \cap B)^c = A^c \cup B^c$ . Secondly,

$$\mathcal{N}(A)^\perp = \mathcal{C}(A^T)$$

and by distributive law, we have

$$\left(\mathcal{C}(A^T) \oplus \Omega^\perp\right) \cap \Omega = \mathcal{C}(A^T) \cap \Omega$$

But the fact  $\Omega$  is a vector space implies

$$\Omega = \mathcal{C}(P_\Omega)$$

Therefore, we have

$$\mathcal{W}^\perp \cap \Omega = \mathcal{C}(A^T) \cap \mathcal{C}(P_\Omega) = \mathcal{C}(P_\Omega A^T)$$

□

### 3.5 Constrained least square estimation

We can apply the results from previous section to obtain constrained least square estimation for  $\beta$  in the following problem:

$$\begin{aligned} \min_{\beta} \quad & \|y - X\beta\|_2^2 \\ \text{s.t.} \quad & A\beta = c \end{aligned} \tag{3.15}$$

where  $A$  is a  $q \times p$  matrix with  $\text{rank}(A) = q$ .

Let  $A\beta_0 = c$  where  $\beta_0$  is a particular solution. Reparametrize the model as follows:

$$\begin{aligned} \mathbb{E}(y - X\beta_0) &= X\beta - X\beta_0 = X\gamma, \gamma = \beta - \beta_0 \\ \mathbb{E}\tilde{y} &= X\gamma = \theta \in \mathcal{C}(X) = \Omega \end{aligned}$$

where  $\tilde{y} = y - X\beta_0$ . Hence,

$$A(X^T X)^{-1} X^T \theta = A(X^T X)^{-1} X^T X \gamma = A\beta - A\beta_0 = 0$$

Thus,

$$\theta \in \mathcal{W} := \mathcal{N}(A(X^T X)^{-1} X^T) \cap \Omega \tag{3.16}$$

By previous section, the projection operator onto  $\mathcal{W}$  is indeed

$$P_{\mathcal{W}} = P_{\Omega} - P_{\mathcal{W}^{\perp} \cap \Omega}$$

But the first term is  $P_{\Omega} = X(X^T X)^{-1} X^T$  and the second term is

$$P_{\mathcal{W}^{\perp} \cap \Omega} = X(X^T X)^{-1} A^T (A(X^T X)^{-1} A^T)^{-1} A(X^T X)^{-1} X^T$$

Let  $\hat{\beta}_H$  be the constrained LSE, then we have

$$\hat{\theta} = P_{\mathcal{W}} \tilde{y} = X(\hat{\beta}_H - \beta_0)$$

Recall the ordinary LSE  $\hat{\beta} = (X^T X)^{-1} X^T y$  and  $c = A\beta_0$ , thus,

$$\begin{aligned} X\hat{\beta}_H &= \hat{\theta} + X\beta_0 \\ &= (P_{\Omega} - P_{\mathcal{W}^{\perp} \cap \Omega}) \tilde{y} + X\beta_0 \\ &= X\hat{\beta} - X(X^T X)^{-1} A^T (A(X^T X)^{-1} A^T)^{-1} (A\hat{\beta} - c) \end{aligned}$$

Since  $X$  is of full rank  $p$ , pre-multiply both sides by  $(X^T X)^{-1} X^T$ :

$$\hat{\beta}_H = \hat{\beta} - (X^T X)^{-1} A^T (A(X^T X)^{-1} A^T)^{-1} (A\hat{\beta} - c) \quad (3.17)$$

### 3.6 When $X$ has less than full column rank

Suppose now the model is

$$\begin{aligned} \underbrace{y}_{n \times 1} &= \underbrace{X}_{n \times p} \underbrace{\beta}_{p \times 1} + e \\ \text{rank}(X) &= r < p \\ \mathbb{E}e &= 0 \\ \text{Cov}(e) &= \sigma^2 I_n \end{aligned}$$

The normal equation  $X^T X \beta = X^T y$  always has a solution

$$\hat{\beta} = (X^T X)^- X^T y$$

Let  $\theta = \mathbb{E}y = X\beta$ , then  $\hat{\theta} = X\hat{\beta} = X(X^T X)^- X^T y$  and  $\hat{\theta}$  is the orthogonal projection of  $y$  onto the space  $\Omega = \mathcal{C}(X)$ . Further, the orthogonal projection operator  $X(X^T X)^- X^T$  is unique. The question is, can I find an unbiased linear estimator of  $\beta$ , i.e., is there a matrix  $C$ , such that

$$\mathbb{E}Cy = \beta.$$

Unfortunately, the answer is "NOT possible".

$$\text{LHS} = \mathbb{E}Cy = CX\beta = \beta = \text{RHS}$$

The above holds for arbitrary  $\beta$ , hence,

$$CX = I$$

But  $X$  is rank-deficient:  $\text{rank}(X) = r < p$ , thus a contradiction.

**Definition 3.1.** Let  $c$  be a  $p$ -dimensional vector. Then  $c^T\beta$  is called estimable if there exists a vector  $a$  such that  $\mathbb{E}a^Ty = c^T\beta$ .

**Theorem 3.3** (Estimability). The parameter  $c^T\beta$  is estimable if and only if

$$c^T = c^T(X^TX)^-X^TX. \quad (3.18)$$

*Proof.* ( $\Leftarrow$ ) The "if" part is trivial:

$$\mathbb{E}c^T(X^TX)^-X^Ty = c^T(X^TX)^-X^TX\beta = c^T\beta$$

In other words,  $a^T = c^T(X^TX)^-X^T$ .

( $\Rightarrow$ ) If  $c^T\beta$  is estimable, then there exists a vector  $a$  s.t.

$$\mathbb{E}a^Ty = a^TX\beta = c^T\beta$$

for any  $\beta$ . Thus  $c = X^Ta$ . Substitute it into RHS of the equality:

$$c^T(X^TX)^-X^TX = a^TX(X^TX)^-X^TX$$

But the vector  $X^Ta$  can be re-written as  $X^TXb$  for some  $b$ ,

$$\begin{aligned} a^TX(X^TX)^-X^TX &= b^TX^TX(X^TX)^-X^TX \\ &= b^TX^TX \\ &= a^TX \\ &= c^T \end{aligned}$$

where the second equality comes from the definition of  $g$ -inverse. □

**Example 3.4** (Fixed effect model). The simplest one-way ANOVA model is:

$$\begin{aligned} y_{ij} &= \alpha + \tau_i + \epsilon_{ij} \\ i &= 1, 2, \dots, K. \\ j &= 1, 2, \dots, n_i. \end{aligned}$$



In matrix form, we have

$$\mathbb{E}y = X\beta = \begin{pmatrix} 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & 0 & 0 & \cdots & 1 \end{pmatrix} \begin{pmatrix} \alpha \\ \tau_1 \\ \cdots \\ \tau_k \end{pmatrix}$$

Suppose we are interested in estimating

$$c^T \beta = c_0 \alpha + c_1 \tau_1 + \cdots + c_K \tau_K$$

By the result 1.2 obtained from Analissa's transformation, we have

$$(X^T X)^- X^T X = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & 0 & 0 & \cdots & 1 \end{pmatrix} \quad (3.19)$$

This is just one of the infinite choices of  $(X^T X)^-$ , but  $c^T (X^T X)^- (X^T X)$  is unique and it is

$$c^T = \begin{pmatrix} c_0 & c_1 & \cdots & c_K \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^K c_i & c_1 & \cdots & c_K \end{pmatrix}$$

In conclusion, if  $c^T \beta$  is estimable, then the vector  $c$  must satisfy the linear constraint:

$$c_0 = \sum_{i=1}^K c_i \quad (3.20)$$

In the 250A final, there is a problem that extends one-way ANOVA to two-way ANOVA fixed effect models (exercise).

### 3.7 Some testing problems

Consider the following model

$$\mathbb{E}y = X\beta, \text{ rank}(X) = p,$$

and we want test

$$H_0 : A\beta = c$$

where  $A$  is  $q \times p$  matrix of full row rank. By previous sections, we have

$$\begin{aligned}\hat{\beta}_H &= \arg \min_{A\beta=c} \|y - X\beta\|_2^2 \\ \hat{\beta} &= \arg \min_{\theta \in \mathcal{C}(X)} \|y - \underbrace{X\beta}_{\theta}\|_2^2\end{aligned}$$

Express  $\hat{\beta}_H$  in terms of  $\hat{\beta}$  (exercise),

$$\|y - X\hat{\beta}_H\|_2^2 = \underbrace{\|y - X\hat{\beta}\|_2^2}_{\in \Omega^\perp} + \underbrace{\|X\hat{\beta} - X\hat{\beta}_H\|_2^2}_{\in \Omega}$$

where  $\Omega = \mathcal{C}(X)$ . For the second term, we have

$$\begin{aligned}\|X\hat{\beta} - X\hat{\beta}_H\|_2^2 &= \text{SSE}_0 - \text{SSE} \\ &= (A\hat{\beta} - A\beta)^T \frac{A(X^T X)^{-1} A^T}{\sigma^2} (A\hat{\beta} - A\beta)\end{aligned}\tag{3.21}$$

But we have

$$A\hat{\beta} - c \sim \mathcal{N}(A\beta - c, \sigma^2 A(X^T X)^{-1} A^T)$$

and hence,

$$3.21 \sim \chi_q^2(\delta^2)\tag{3.22}$$

To construct a test statistic for  $H_0$ , we have

$$\frac{\text{SSE}_0 - \text{SSE}}{\text{SSE}} \frac{n - q}{q} \sim_{H_0} F_{q, n-p}(0)\tag{3.23}$$

**Example 3.5** ( $K$ -phase regression). Suppose the model is:

$$y_{ki} = \alpha_k + \beta_k x_k + \epsilon_{ki}, \quad i = 1, 2, \dots, n_k; \quad k = 1, 2, \dots, K.$$

We want to test the hypothesis

$$H_0 : \tilde{\beta} = \beta_1 = \beta_2 = \dots = \beta_K$$

In other words, whether the slopes of  $K$  regression lines are the same or not.

In matrix form, we have

$$A\beta = \begin{pmatrix} 0 & 0 & \cdots & 0 & 1 & 0 & \cdots & 0 & -1 \\ 0 & 0 & \cdots & 0 & 0 & 1 & \cdots & 0 & -1 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 1 & -1 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \cdots \\ \alpha_K \\ \beta_1 \\ \beta_2 \\ \cdots \\ \beta_K \end{pmatrix}$$

and  $\text{rank}(A) = K - 1$ . Thus, to test the hypothesis, we have to solve for the  $\text{SSE}_H$  with constraint

$$A\beta = 0$$

The test statistic is

$$F = \frac{\text{SSE}_H - \text{SSE}}{\text{SSE}} \frac{N - 2K}{K - 1} \sim_{H_0} F_{K-1, N-2K} \quad (3.24)$$

where

$$\begin{aligned} N &= n_1 + n_2 + \cdots + n_K \\ \text{SSE} &= \sum_i \sum_j (y_{ij} - \alpha_i - \beta_i x_{ij})^2 \\ \text{SSE}_H &= \sum_i \sum_j (y_{ij} - \alpha_i - \tilde{\beta} x_j)^2 \end{aligned}$$

**Example 3.6** ( $K$ -phase regression continued). Suppose now we want to test whether  $K$ -regression lines are the same or not, then the constraint becomes

$$0 = A\beta = \begin{pmatrix} 0 & 0 & \cdots & 0 & 0 & 1 & 0 & \cdots & 0 & -1 \\ 0 & 0 & \cdots & 0 & 0 & 0 & 1 & \cdots & 0 & -1 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 0 & 0 & 0 & 0 & \cdots & 1 & -1 \\ 1 & 0 & \cdots & 0 & -1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & -1 & 0 & 0 & \cdots & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 1 & -1 & 0 & 0 & \cdots & 0 & 0 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \cdots \\ \alpha_K \\ \beta_1 \\ \beta_2 \\ \cdots \\ \beta_K \end{pmatrix}$$

**Example 3.7** ( $K$ -phase regression continued). For a fixed point  $\phi$ , we want to test the hypothesis:

$$H_0 : \alpha_i + \beta_i \phi = \text{constant for } i = 1, 2, \cdots, K.$$

The general case when  $\text{rank}(X) = r < p$  is similar to the full column rank case. So I omit them.

### 3.8 Exercises

1. Prove Gauss-Markov theorem. 3.1
2. Verify inequality 3.10 using the fact that

$$\text{Cov}(\hat{\beta}, (X^T X)^{-1} X^T Z \hat{\gamma}_G) = 0$$

3. Verify lemma 10.
4. ([1] McElroy, 1967) Let  $Y_1, Y_2, \dots, Y_n$  be random variables with common mean  $\theta$  and with dispersion matrix  $\sigma^2 V$ , where  $v_{ii} = 1$  ( $i = 1, 2, \dots, n$ ) and  $v_{ij} = \rho$  ( $0 < \rho < 1; i, j = 1, 2, \dots, n; i \neq j$ ). Find the generalized least square estimator of  $\theta$  and show that it is the same as the ordinary least square estimate. (Hint:  $V^{-1}$  takes the same form as  $V$ ).
5. Show that the matrix

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

is UNIQUE, i.e., it does not depend on the g-inverse (Hint: show that  $P$  is an OP matrix, and then use the fact that OP is unique).

6. (250A Final) In a least square fit using the standard linear model  $\mathbb{E}y = Xb$  with 8 points, the fitted values for the second and third data points are given by the equations

$$\hat{y}_2 = -0.242y_1 + 0.374y_2 + 0.329y_3 + 0.182y_4 + 0.009y_5 - 0.115y_6 - 0.113y_7 + 0.091y_8$$

$$\hat{y}_3 = -0.061y_1 + 0.329y_2 + 0.418y_3 + 0.312y_4 + 0.117y_5 - 0.061y_6 - 0.115y_7 + 0.061y_8$$

The error sum of squares is 400 on 4 degrees of freedom. Assume that all errors are independent, normally distributed errors with common variance  $\sigma^2$ .

- (a) Provide where possible, estimates for the variances and covariances of the residuals  $e_1, e_2$  and  $e_3$  at the first three data points.
  - (b) Is it plausible that there are 2 other rows in the  $X$  matrix that are the same as the second row? Justify.
7. (250A Final) Is it true that for any orthogonal projection matrix  $P$  with elements  $p_{ij}$ ,
    - (a)  $p_{ij} \leq (p_{ii} p_{jj})^{1/2}$ ?
    - (b)  $(1 - p_{ii})(1 - p_{jj}) - p_{ij}^2 \geq 0$ ?

Provide justifications for your answers.

8. Show that if  $X$  has full rank,

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

where  $\hat{\beta}$  is the least square estimate. Such relation is usual when we are deriving Bayes estimator in linear models in chapter 7.

9. (250A Final) Suppose  $Y$  is the  $n \times 1$  vector of observations and the rank of the  $n \times p$  matrix  $X$  is  $r < p$  in the below model.

$$Y = X\beta + \epsilon, \epsilon \sim \mathcal{N}_n(0, \sigma^2 I)$$

(i) Answer the following questions True (T) of False (F) with justifications:

- There is no linear unbiased estimator for  $\beta$ .
- If  $a^T \beta$  and  $c^T \beta$  are linear estimable functions, then  $(a + c)^T \beta$  is also estimable.
- The sum of two linear non-estimable functions can be estimable.

(ii) Consider the balanced two-way ANOVA with interactions:

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$$

$i = 1, 2, \dots, a; j = 1, 2, \dots, b$  and  $k = 1, 2, \dots, n$ .

- Write down the model in matrix form making clear the explicit form of the design matrix  $X$ .
- Identify all linear estimable functions for the model.
- Does your answer in the above depend on whether the numbers of observations in the cells are equal? Justify your answer.

10. Let

$$Y_i = \beta_0 + \beta_1(x_{i1} - \bar{x}_1) + \beta_2(x_{i2} - \bar{x}_2) + \epsilon_i, i = 1, 2, \dots, n$$

where  $\bar{x}_j = \sum_{i=1}^n x_{ij}/n$ ,  $\mathbb{E}\epsilon = 0$ , and  $Var(\epsilon) = \sigma^2 I_n$ . If  $\hat{\beta}_1$  is the least square estimate of  $\beta_1$ , show that

$$Var(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2 (1 - r^2)}$$

where  $r$  is the correlation coefficient of the  $n$  pairs  $(x_{i1}, x_{i2})$ .

11. Apply the spectral decomposition theorem to find the maximum likelihood estimation of  $\mu$  and  $\Sigma$  if  $y_1, \dots, y_n \sim \mathcal{N}_p(\mu, \Sigma)$ .

**Hint:** the log-likelihood is

$$C - \frac{n}{2} \log \det \Sigma - \frac{1}{2} \sum_{i=1}^n (y_i - \bar{y})^T \Sigma^{-1} (y_i - \bar{y}) - \frac{n}{2} (\bar{y} - \mu)^T \Sigma^{-1} (\bar{y} - \mu)$$

where  $\bar{y} = \sum_{i=1}^n y_i/n$ . But the last term is non-negative and is 0 iff  $\bar{y} = \mu$ . Hence,

$$\hat{\mu}_{MLE} = \bar{y}$$

Next, let  $A = \sum_{i=1}^n (y_i - \bar{y})(y_i - \bar{y})^T$  and the log-likelihood with  $\hat{\mu}_{MLE}$  becomes

$$C - \frac{n}{2} \log \det \Sigma - \frac{n}{2} \text{Tr}(A^{1/2} \Sigma^{-1} A^{1/2})$$

Since  $A$  is positive definite almost surely (a theorem), the problems becomes

$$\begin{aligned} & \max [\log \det TDT^T - \text{Tr} TDT^T] \\ & \text{s.t. } A^{1/2} \Sigma^{-1} A^{1/2} = TDT^T \text{ is the spectral decomposition.} \end{aligned} \quad (3.25)$$

But this is equivalent to

$$\max \left( \sum_{i=1}^p \log \lambda_i - \sum_{i=1}^p \lambda_i \right) \quad (3.26)$$

where  $\lambda_i$ 's are eigen-values of  $A^{1/2} \Sigma^{-1} A^{1/2}$ . After some calculus, we can get

$$\hat{\Sigma}_{MLE} = A = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(y_i - \bar{y})$$

12. (A useful property of OP) Show that all off-diagonal entries in an orthogonal projection matrix are between  $-0.5$  and  $0.5$ .

**Hint:** use the facts that  $\sum_j p_{ij} p_{ji} = p_{ii}$  and  $p_{ij} = p_{ji}$ .

13. (Variance of CLSE) Let  $\hat{\beta}_H$  be the constrained least square estimate, prove that

$$\text{Var}(\hat{\beta}_H) = \sigma^2 ((X^T X)^{-1} - (X^T X)^{-1} A^T (A(X^T X)^{-1} A^T)^{-1} A(X^T X)^{-1})$$

Hence deduce that

$$\text{Var}(\hat{\beta}_{Hj}) \leq \text{Var}(\hat{\beta}_j)$$

14. Let  $\hat{Y}_H$  be the constrained prediction of  $Y$ , i.e.,  $\hat{Y}_H = X\hat{\beta}_H$ . Prove that

$$\|Y - \hat{Y}_H\|_2^2 = \|Y - \hat{Y}\|_2^2 + \|\hat{Y} - \hat{Y}_H\|_2^2.$$

15. (Largangian multipler) If  $\hat{\lambda}_H$  is the least square estimate of the Lagrange multiplier associated with the constraints  $A\beta = c$ , show that

$$\text{RSS}_H - \text{RSS} = \sigma^2 \hat{\lambda}_H^T \left( \text{Var}(\hat{\lambda}_H) \right)^{-1} \hat{\lambda}_H$$

This idea is used to construct Lagrange multiplier tests.

16. (Invariance property of  $F$ -statistic) Let  $Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_{p-1} x_{i,p-1} + \epsilon_i, i = 1, 2, \dots, n$ , where the  $\epsilon_i$  are independent  $\mathcal{N}(0, \sigma^2)$ . Prove that the  $F$ -statistic for testing the hypothesis  $H : \beta_q = \beta_{q+1} = \cdots = \beta_{p-1} = 0$  ( $0 < q \leq p-1$ ) is unchanged if a constant,  $c$ , say, is subtracted from each  $Y_i$ .
17. Given  $Y = \theta + \epsilon$ , where  $\epsilon \sim \mathcal{N}_4(0, \sigma^2 I_4)$  and  $\theta_1 + \theta_2 + \theta_3 + \theta_4 = 0$ , show that the  $F$ -statistic for testing  $H : \theta_1 = \theta_3$  is

$$\frac{2(Y_1 - Y_3)^2}{(Y_1 + Y_2 + Y_3 + Y_4)^2}.$$

## 4 Multiple and Partial Correlation Coefficient

In this section, I study the basic properties of LS coefficient and multiple and partial correlation coefficient (MCC & PCC).

### 4.1 Multiple correlation coefficient (MCC)

The multiple correlation coefficient is denoted as  $r$  and is defined as

$$r = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{(\sum_{i=1}^n (y_i - \bar{y})^2)(\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2)}}$$

**Theorem 4.1** (Relation between  $r$  and  $R^2$ ). If  $\mathbf{1} \in \mathcal{C}(X)$ , then  $r^2 = R^2$ .

*Proof.* If  $\mathbf{1} \in \mathcal{C}(X)$ , then  $\hat{y} = Py \Rightarrow \mathbf{1}^T \hat{y} = \mathbf{1}^T Py = (P\mathbf{1})^T y = \mathbf{1}^T y \Rightarrow \bar{\hat{y}} = \bar{y}$ . Thus,

$$\sum_{i=1}^n \bar{y}(\hat{y}_i - \bar{\hat{y}}) = y^T P_1 (Py - P_1 y) = 0$$

Thus I re-write  $r = \frac{y^T (P - P_1)y}{\sqrt{y^T Q_1 y y^T (P - P_1)y}}$ . Therefore,  $R^2$  is

$$R^2 = \frac{\text{SSR}}{\text{TSS}} = \frac{y^T (P - P_1)y}{y^T Q_1 y} = r^2$$

□

**Example 4.1.** We want to test  $\beta_1 = \dots = \beta_{p-1} = 0$  in a LM. This boils down to a  $F$ -statistic:

$$\begin{aligned} F &= \frac{\text{SSE}_0 - \text{SSE}}{\text{SSE}} \frac{n-p}{p-1} \\ &= \frac{\text{TSS} - \text{SSR}_0 - (\text{TSS} - \text{SSR})}{\text{TSS} - \text{SSR}} \frac{n-p}{p-1} \\ &= \frac{1 - R_0^2 - (1 - R^2)}{1 - R^2} \frac{n-p}{p-1} \\ &= \frac{R^2 - R_0^2}{1 - R^2} \frac{n-p}{p-1} \\ &= \frac{R^2}{1 - R^2} \frac{n-p}{p-1} \end{aligned}$$



The last equality is due to  $R_0^2 = \frac{y^T(P_1 - P_1)y}{\text{TSS}} = 0$ . Thus, test statistic for  $H_0$  is

$$F = \frac{R^2}{1 - R^2} \frac{n - p}{p - 1} \quad (4.1)$$

$$R^2 = \frac{(p - 1)F}{(n - p) + (p - 1)F}$$

Note (exercise)  $X \sim F_{d_1, d_2} \Rightarrow \frac{d_1 X / d_2}{1 + d_1 X / d_2} \sim B(\frac{d_1}{2}, \frac{d_2}{2})$ . So

$$R^2 \sim B(\frac{p - 1}{2}, \frac{n - p}{2}), \mathbb{E}(R^2) = \frac{p - 1}{n - 1}$$

## 4.2 Geometry of LSE

Let  $\mathcal{V}_{-k} = \mathcal{L}(x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_p)$  and  $\theta = \mathbb{E}(y) = X\beta = \sum_{i=1}^p \beta_i \mathbf{x}_i$  and for  $k = 1, \dots, p$ ,

$$\hat{\mathbf{x}}_k = P_{\mathcal{V}_{-k}}(\mathbf{x}_k) := P_{-k} \mathbf{x}_k$$

$$\mathbf{x}_k^\perp = \mathbf{x}_k - \hat{\mathbf{x}}_k = Q_{\mathcal{V}_{-k}}(\mathbf{x}_k) := Q_{-k} \mathbf{x}_k$$

and  $\mathbf{x}_k^\perp$  provides information other than these provided by  $\{x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_p\} \subset \mathcal{V}_{-k}$ .

**Example 4.2** (Geometry of LSE). Observe that for  $j = 1, \dots, p$ ,

$$\langle \theta, \mathbf{x}_j^\perp \rangle = \langle \sum_{i=1}^p \beta_i \mathbf{x}_i, \mathbf{x}_j^\perp \rangle = \beta_j \|\mathbf{x}_j^\perp\|_2^2$$

Replace  $\theta$  with  $y$  and  $\beta_i$  with  $\hat{\beta}_i$ , then for  $i, j = 1, \dots, p$ , we have

$$\hat{\beta}_i = \frac{\langle y, \mathbf{x}_i^\perp \rangle}{\|\mathbf{x}_i^\perp\|_2^2}$$

$$= \frac{y^T Q_{-i} \mathbf{x}_i}{\mathbf{x}_i^T Q_{-i} \mathbf{x}_i} \quad (4.2)$$

$$\text{Var}(\hat{\beta}_i) = \frac{\sigma^2}{\|\mathbf{x}_i^\perp\|_2^2}$$

$$= \frac{\sigma^2}{\mathbf{x}_i^T Q_{-i} \mathbf{x}_i} \quad (4.3)$$

$$\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = \frac{\sigma^2 \langle \mathbf{x}_i^\perp, \mathbf{x}_j^\perp \rangle}{\|\mathbf{x}_i^\perp\|_2^2 \|\mathbf{x}_j^\perp\|_2^2}$$

$$= \cos(\mathbf{x}_i^\perp, \mathbf{x}_j^\perp) \frac{\sigma^2}{\|x_i\|_2 \|x_j\|_2} \quad (4.4)$$

$$\sigma^2 (X^T X)_{ij}^{-1} = \cos(\mathbf{x}_i^\perp, \mathbf{x}_j^\perp) \frac{\sigma^2}{\|x_i\|_2 \|x_j\|_2} \quad (4.5)$$

**Lemma 12** (F-test). Define

$$R^2 = \frac{SSR}{TSS}, R_{k-1}^2 = \frac{SSR(\mathbf{x}_1, \dots, \mathbf{x}_{k-1})}{TSS}, R_k^2 = \frac{SSR(\mathbf{x}_1, \dots, \mathbf{x}_k)}{TSS}$$

Then we have

$$\frac{R_k^2 - R_{k-1}^2}{1 - R_{k-1}^2} = \frac{t^2}{n - k + t^2} \quad (4.6)$$

where  $t$  is the test statistic for  $\beta_k = 0$ .

### 4.3 Partial correlation coefficient (PCC)

**Definition 4.1** (PCC). Let  $\mathbf{v}_1, \mathbf{v}_2, \mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$ , then the partial correlation coefficient (PCC) of  $\mathbf{v}_1$  and  $\mathbf{v}_2$  with the linear effects of  $\mathbf{x}_1, \dots, \mathbf{x}_k$  removed is

$$r_{\mathbf{v}_1 \mathbf{v}_2 \cdot \mathbf{x}_1 \dots \mathbf{x}_k} = \frac{\langle \mathbf{v}_1^\perp, \mathbf{v}_2^\perp \rangle}{\|\mathbf{v}_1^\perp\|_2 \|\mathbf{v}_2^\perp\|_2}$$

where  $\mathbf{v}_i^\perp = Q \mathbf{v}_i, i = 1, 2$ . When  $k = 1, \mathbf{x}_1 = \mathbf{1}$ , then  $r_{\mathbf{v}_1 \mathbf{v}_2 \cdot \mathbf{x}_1}$  is just Pearson correlation.

Now we express PCC of order  $k$  in terms of order  $k - 1$ . For example, if  $k = 4$ , then

$$r_{14 \cdot 23} = \frac{r_{14} - r_{13 \cdot 2} r_{34 \cdot 2}}{\sqrt{1 - r_{13 \cdot 2}^2} \sqrt{1 - r_{34 \cdot 2}^2}}$$

Consider the general case in the definition and let  $\mathcal{V}_J = \mathcal{L}(x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_k)$  and  $\mathcal{V}_I = \mathcal{L}(x_1, \dots, x_k)$ . Define  $P_J = P_{\mathcal{V}_J}, Q_J = I - P_J, P_I = P_{\mathcal{V}_I}, Q_I = I - P_I$  and

$$\hat{v}_i := P_J \mathbf{v}_i, v_i^\perp := Q_J \mathbf{v}_i, \hat{x}_i := P_J \mathbf{x}_i, x_i^\perp := Q_J \mathbf{x}_i$$

By 1.11 we have  $P_I = P_J + \frac{Q_J \mathbf{x}_j \mathbf{x}_j^T Q_J}{\mathbf{x}_j^T Q_J \mathbf{x}_j}$ . Thus,

$$\begin{aligned} w_i &= \mathbf{v}_i - P_I \mathbf{v}_i \\ &= Q_J \mathbf{v}_i - \frac{Q_J \mathbf{x}_j \mathbf{x}_j^T Q_J}{\mathbf{x}_j^T Q_J \mathbf{x}_j} \mathbf{v}_i \\ &= v_i^\perp - \frac{\langle x_j^\perp, \mathbf{v}_i^\perp \rangle}{\|x_j^\perp\|_2^2} x_j^\perp \end{aligned}$$

Thus, for  $i = 1, 2$ ,

$$\|w_i\|_2^2 = \|v_i^\perp\|_2^2 - \frac{\langle x_j^\perp, \mathbf{v}_i^\perp \rangle^2}{\|x_j^\perp\|_2^2} = \|v_i^\perp\|_2^2 (1 - r_{ij \cdot J}^2)$$

Where  $r_{ij \cdot J} = r_{\mathbf{v}_1 \mathbf{v}_2 \cdot \mathbf{x}_1 \dots \mathbf{x}_{j-1} \mathbf{x}_{j+1} \dots \mathbf{x}_k}$ . Besides,

$$\langle w_1, w_2 \rangle = \|v_1^\perp\|_2 \|v_2^\perp\|_2 (r_{12 \cdot J} - r_{13 \cdot J} r_{23 \cdot J})$$

To sum up, we have shown

$$r_{12 \cdot I} = \frac{\langle w_1, w_2 \rangle}{\|w_1\|_2 \|w_2\|_2} = \frac{r_{12 \cdot J} - r_{13 \cdot J} r_{23 \cdot J}}{\sqrt{1 - r_{1j \cdot J}^2} \sqrt{1 - r_{2j \cdot J}^2}} \quad (4.7)$$

## 4.4 Exercises

1. If  $X$  has an intercept (all ones), prove that

$$\sum_{i=1}^n (Y_i - \hat{Y}_i) = 0$$

Hint: consider the first column of  $X$ .

2. If  $F \sim F_{a,b}$ , then show that

(a)  $\frac{aF/b}{1+aF/b} \sim B(a/2, b/2)$ .

(b) The expectation of this Beta RV is  $\frac{a}{a+b}$ .

3. ([1], p113) Consider a general full rank linear regression model. Show that  $R^2$  and  $F$  for testing  $H_0 : \beta_j = 0, j \neq 0$  are independent of the units in which  $y_i$  and  $x_{ij}$  are measured.
4. Let  $y = \beta_0 + \sum_{i=1}^k \beta_i x_i$ . Show  $r_{yx_k \cdot x_1 \dots x_{k-1}}$  is a function of the test statistic for  $H_0 : \beta_k = 0$ .

## 5 Violation of Assumptions and Remedies

This section talks about outlier detection, consequences of under-fitting and over-fitting, transformation and collinearity.

### 5.1 Outlier detection

Consider a general linear model:

$$\begin{aligned} y &= X\beta + \epsilon \\ \mathbb{E}y &= X\beta \\ \hat{y} &= X\hat{\beta} = Py \\ e &= y - \hat{y} = Qy \end{aligned}$$

We have  $\mathbb{E}e = 0$  and  $Var(e) = \sigma^2 Q$ ,  $Cov(e, \hat{y}) = Cov(Qy, Py) = 0$ .

**Definition 5.1** (Internally and externally studentized residual). The internally studentized residual is defined as

$$r_i = \frac{e_i}{s\sqrt{(1 - h_{ii})}}$$

and the externally studentized residual is defined as

$$t_i = \frac{e_i}{s_{(i)}\sqrt{1 - h_{ii}}}$$

where  $h_{ii}$  is the  $i^{th}$  diagonal element of  $P$  and  $s, s_{(i)}$  are the usual estimates of  $\sigma^2$  with and without the  $i^{th}$  case.

Recall that we have

$$\hat{\beta} - \hat{\beta}_{(i)} = \frac{(X^T X)^{-1} x_i e_i}{1 - h_{ii}}, \quad e_i = y_i - x_i^T \hat{\beta}$$

where  $x_i$  is the  $i^{th}$  row of  $X$ . We can use this formula to show the relationship between  $s_{(i)}^2$  and  $s^2$ . Besides, we have two useful results:

$$\begin{aligned} t_i^2 &=_d \frac{B}{1 - B}(n - p - 1) \\ r_i^2 &=_d B(n - p) \end{aligned} \tag{5.1}$$

where  $B \sim Beta(1/2, (n - p - 1)/2)$ . For proof, see exercises.

### 5.2 Under-fitting and over-fitting

In this subsection, we study the linear model when the working model is over and under fitted.

**Example 5.1** (Underfitting). Suppose the working model is (full rank)

$$\mathbb{E}(y) = X\beta$$

and the true model is

$$\mathbb{E}(y) = X\beta + Z\gamma$$

Let  $\hat{\beta} = (X^T X)^{-1} X^T y$  be the usual LSE and  $Cov(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$ . Then,

$$\mathbb{E}\hat{\beta} = \beta + (X^T X)^{-1} X^T Z\gamma$$

Thus,  $\hat{\beta}$  is unbiased iff  $X^T Z = 0$ . Next, SSE is

$$s^2 = \frac{y^T Q y}{n - p}$$

Thus,

$$\begin{aligned} \mathbb{E}s^2 &= \frac{1}{n - p} ((\mathbb{E}y)^T Q \mathbb{E}y + Tr(Q Var(y))) \\ &= \sigma^2 + \frac{1}{n - p} \gamma^T Z^T Q Z \gamma \end{aligned}$$

So  $s^2$  overestimates  $\sigma^2$  unless  $QZ = 0$ . In addition, let  $\hat{y} = X\hat{\beta}$ ,

$$\mathbb{E}\hat{y} = X\beta + P_X Z\gamma$$

Finally, let  $e = y - \hat{y}$ , then

$$\begin{aligned} \mathbb{E}e &= \mathbb{E}y - \mathbb{E}X\hat{\beta} = Q_X Z\gamma \\ Var(e) &= \sigma^2 Q \end{aligned}$$

So  $Var(e)$  stays the same. Next, let  $\tilde{\beta}$  be the LSE under the true model, then

$$\begin{aligned} \tilde{y} &= (X \ Z)\tilde{\beta} \\ &= Py - QZ(Z^T QZ)^{-1} Z^T Qy \end{aligned}$$

Thus,  $Var(\tilde{y}) - Var(\hat{y})$  is positive semi-definite.

**Example 5.2** (Overfitting). The true model is  $\mathbb{E}y = X_1\beta_1$  and the working model is  $\mathbb{E}y = X_1\beta_1 + X_2\beta_2$ . Let  $X = (X_1 \ X_2)$ , then the expectation of LSE is

$$\mathbb{E}\hat{\beta} = (X^T X)^{-1} X^T (X_1 \ X_2) (\beta_1^T \ 0)^T$$

So  $\hat{\beta}_1$  is unbiased.

**Example 5.3** (Mis-specification of covariance). The working model is  $Cov(\epsilon) = \sigma^2 I$  while the true

model is  $Cov(\epsilon) = \sigma^2 V$ . Then  $\hat{\beta}$  is still unbiased but

$$Cov(\hat{\beta}) = \sigma^2 (X^T X)^{-1} X^T V X (X^T X)^{-1}$$

But is  $\hat{\sigma}^2 = \frac{y^T Q y}{n-p}$  biased for  $\sigma^2$ ?

**Theorem 5.1** (Bounding the ratio). If  $\mu_1 \leq \mu_2 \leq \dots \leq \mu_n$  are eigenvalues of  $V$ , then

$$\frac{1}{n-p} \sum_{i=1}^{n-p} \mu_i \leq \mathbb{E} \frac{\hat{\sigma}^2}{\sigma^2} \leq \frac{1}{n-p} \sum_{i=n-p+1}^n \mu_i$$

To prove theorem 5.1, we need the following lemma.

**Lemma 13.** Let  $H^T = H = H^2$  and  $rank(H) = r$ . If  $A^T = A$ , then

$$\sum_{i=1}^r \lambda_i \leq Tr(HA) \leq \sum_{i=n-r+1}^n \lambda_i$$

where  $\lambda_1 \leq \dots \leq \lambda_n$  are eigenvalues of  $A$ .

*Proof.* By spectral decomposition theorem,  $\exists P$  and  $P^T P = P^T P = I$  such that

$$P^T H P = \begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix} := D_r$$

Note  $D_r^2 = D_r$ . Thus,

$$\begin{aligned} Tr(HA) &= Tr(P D_r P^T A) \\ &= Tr((P D_r)^T (A P D_r)) \\ &= \sum_{i=1}^r p_i^T A p_i \end{aligned}$$

where  $p_i$  is the  $i^{th}$  column of  $A$ . Then the lemma follows from theorem 1.9. □

## 5.3 Transformation

Consider transform  $x$  by  $g(x)$ , then

$$\mathbb{E}(y) = \alpha + \beta^T g(x)$$

Consider simple transformations like

$$g_\lambda(x) = \begin{cases} \ln x, & \lambda = 0 \\ \frac{x^\lambda - 1}{\lambda}, & \lambda \neq 0 \end{cases}$$

In practice, run the regression with

$$\lambda = \{-2, -1, -\frac{1}{2}, -\frac{1}{3}, 0, \frac{1}{3}, \frac{1}{2}, 1, 2\}$$

Do the same procedure on the response ( $y > 0$ ) and ad-hoc. Alternatively, use variance stabilizing transformation so that  $Var(g(y))$  is a constant. By Taylor's expansion,

$$g(y) = g(\mu) + g'(\mu)(y - \mu) + o_p(y - \mu)$$

$$Var(g(y)) \approx Var(y)[g'(\mu)]^2 = \sigma^2[g'(\mu)]^2$$

Thus, to make  $Var(g(y)) = c$  a constant,  $g'(\mu) = \sqrt{\frac{c^2}{\sigma^2}}$ , or,

$$g(\mu) = \int \frac{c}{\sqrt{Var(y)}} d\mu$$

In case of Poisson,  $\mu = Var(y)$ ,  $c = 1$ ,

$$g(y) = 2\sqrt{y}$$

## 5.4 Collinearity

Recall that in a classical linear model, we have least square estimate  $\hat{\beta}^T = (\hat{\beta}_1, \dots, \hat{\beta}_p)$  and

$$Var(\hat{\beta}_k) = \frac{\sigma^2}{\|x_k^\perp\|_2^2}$$

where  $x_k^\perp = x_k - P_{-k}x_k$  and  $P_{-k}$  is the projection operator without the  $k^{th}$  column  $x_k$ . If  $x_k$  can be approximated by linear combinations of the other  $x$ 's (collinearity), then  $\hat{\beta}_k$  will be estimated unreliably. How to detect collinearity?

**Definition 5.2** (Variance inflation factor). Let  $R_j^2$  be the usual  $R^2$  obtained when regressing  $x_j$  on the other covariates, say  $X_{-j}$ . Then variance inflation factor is defined as

$$VIF_j = \frac{1}{1 - R_j^2}$$

If  $VIF_j$  is large, say  $> 10$ , then  $x_j$  is almost linearly related with the rest of  $x$ 's. What to do?

1. Omit the variable or variables with  $VIF > 10$ .
2.  $\mathbb{E}y = X\beta = XR^{-1}R\beta$  where  $R$  is an upper triangular non-singular matrix and is obtained from the QR decomposition of  $X$ :

$$X = QR = \begin{pmatrix} q_1 & q_2 & \cdots & q_p \end{pmatrix} R$$

$$q_i^T q_j = \delta_{ij}, \quad q_1 = \frac{x_1}{\|x_1\|_2}$$

Let  $\gamma = R\beta$  then

$$\hat{\gamma} = \widehat{R\beta} = (Q^T Q)^{-1} Q^T y = Q^T y = (X R^{-1})^T y$$

The regression sum of squares is

$$\text{SSReg} = y^T P_{C(Q)} y = \sum_{i=1}^p (q_i^T y)^2$$

where  $q_i$  is the  $i^{\text{th}}$  column of  $Q$ .

## 5.5 Exercises

1. (Mean shift model) We want to test whether the  $i^{\text{th}}$  case has outlying x-values using the mean-shift mode given by

$$\mathbb{E}y = X\beta + \Theta\phi_i$$

where  $\phi_i$  is the vector with all its components equal to 0 except for the  $i^{\text{th}}$  element, which is equal to 1. Derive a test whether the  $i^{\text{th}}$  case is outlying and show that this test statistic is the  $i^{\text{th}}$  externally studentized statistic and its square has a  $F_{1,n-p-1}$  distribution.

2. ([1], p270) Prove equations 5.1.
3. ([1], p268 and p311) Show that

$$t_i^2 = \frac{(n-p-1)r_i^2}{n-p-r_i^2}$$

4. (Cook's distance) The Cook's distance of the  $i^{\text{th}}$  case is

$$\text{Cook}_i = \frac{(\hat{y} - \hat{y}_{(i)})^T (\hat{y} - \hat{y}_{(i)})}{ps^2}$$

where  $\hat{y}_{(i)} = X\beta_{(i)}$  and  $s^2 = \frac{\text{SSE}}{n-p}$ . Express it in terms of  $p$ ,  $h_{ii}$  and  $r_i$ .

5. (Added variable plot) Suppose  $\mathbb{E}y = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$  and the standard assumptions hold. An added variable plot is constructed by first regressing  $y$  on  $x_1, \dots, x_k$  and regressing  $x_k$  on  $x_1, \dots, x_{k-1}$ , and then regressing the first set of residuals  $e_1$  on the second set of residuals  $e_2$ .
  - (a) What are the fitted regression coefficients when you regress  $e_1$  on  $e_2$  and relate your answers to those fitted coefficients from regressing  $y$  on  $x_1, \dots, x_k$ .
  - (b) What is correlation  $e_1$  and  $e_2$  and describes its relationship to the partial correlation between  $y$  and  $x_k$  controlling for  $x_1, x_2, \dots, x_{k-1}$ .



6. Prove theorem 5.1.

7. (A new metric for detecting outlier) To measure the influence of the  $i^{th}$  observation on the  $i^{th}$  predicted value, it was proposed that a distance measure based on a scale difference  $\hat{y}_i - \hat{y}_{(i)}$  be used:

$$\frac{|\hat{y}_i - \hat{y}_{(i)}|}{s_{(i)} h_{ii}^{1/2}}$$

Here the usual notation holds, eg.  $h_{ii}$  is the leverage for the  $i^{th}$  case, etc.

(a) Show that this statistic can be written as

$$|t_i| \sqrt{\frac{h_{ii}}{1 - h_{ii}}}$$

where  $t_i$  is the externally studentized residual for the  $i^{th}$  case.

(b) Use this statistic and suggest a rule to assess whether the  $i^{th}$  observation is influential.

8. (VIF and LSE) Suppose we regress  $Y$  on  $p$  covariates  $x_1, \dots, x_p$  and assume further that the model has an intercept and all covariates are centered and scaled so that each has sample mean 0 and sample variance equal to unity. Prove or disprove that under standard assumptions, the variance of the LSE for the coefficient of  $j^{th}$  covariate is

$$\frac{\sigma^2}{1 - R_j^2}$$

where  $R_j$  is the sample multiple correlation coefficient obtained by regressing  $x_j$  on the rest of the other  $x_i$ 's. Why is this result useful?

9. (Final exam, covariance ratio) Consider the model  $y = X\beta + \epsilon$ ,  $X$  is a  $n \times p$  matrix of full column rank and  $\mathbb{E}\epsilon = 0$ ,  $Var(\epsilon) = \sigma^2 I_n$ . Let  $b$  be the LSE of  $\beta$  and let  $VR_i$  be the ratio of the determinants of the estimated covariance matrices of  $b$  with and without the  $i^{th}$  case. Express this statistic only in terms of  $n, p, h_{ii}$  and  $r_i$ , the  $i^{th}$  internally studentized residue.

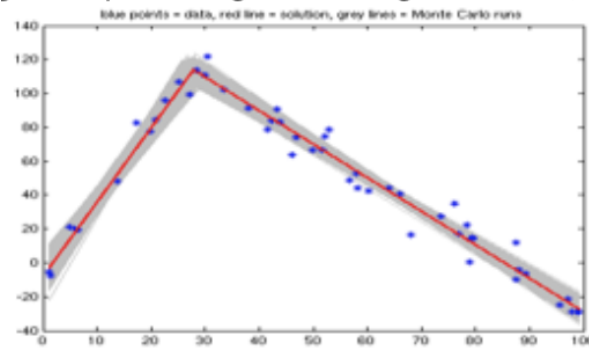
## 6 Hypothesis Testing and Simultaneous Inference

In this section, topics include Fieller's theorem, lack of fitness test, Scheffe's method, Tukey's q. Besides, various tests of homogeneity of variance are discussed, including Hartley's test, Levene's test, Brown-Forsythe's test and Cook's likelihood ratio test.

### 6.1 Fieller's theorem

Consider a two phase regression problem:

Figure 1: Two phase regression.



$$\mathbb{E}y = \begin{cases} \alpha_1 + \beta_1 x, & x \leq \gamma \\ \alpha_2 + \beta_2 x, & x \geq \gamma \end{cases}$$

The continuity assumption

$$\alpha_1 + \beta_1 \gamma = \alpha_2 + \beta_2 \gamma$$

generates the point estimation of  $\gamma$ :

$$\hat{\gamma} = \frac{\hat{\alpha}_1 - \hat{\alpha}_2}{\hat{\beta}_2 - \hat{\beta}_1}$$

The question is: can we construct a confidence interval for  $\gamma$  based on  $\hat{\gamma}$ ? Fieller gave the most complete discussion in 1954.

**Theorem 6.1** (Fieller's theorem). Suppose  $\begin{pmatrix} X \\ Y \end{pmatrix}$  is the mean of a random sample of bi-variate normal random variables with distribution  $\mathcal{N}\left(\begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix}, \sigma^2 \begin{bmatrix} v_{11} & v_{12} \\ v_{21} & v_{22} \end{bmatrix}\right)$ , and  $s^2$  is an unbiased estimator of  $\sigma^2$  such that

$$\frac{rs^2}{\sigma^2} \sim \chi_r^2$$

Let  $\theta = \frac{\mu_X}{\mu_Y}$  be the parameter of interest. Then a  $(1 - 2\alpha)$  confidence interval  $(\theta_L, \theta_U)$  for  $\theta$  is

$$(\theta_L, \theta_U) = \frac{1}{1 - g} \left[ \frac{X}{Y} - \frac{gv_{12}}{v_{22}} \mp \frac{t_{r,\alpha}s}{Y} \sqrt{v_{11} - 2\frac{X}{Y}v_{12} + \frac{X^2}{Y^2}v_{22} - g \left( v_{11} - \frac{v_{12}^2}{v_{22}} \right)} \right] \quad (6.1)$$

where

$$g = \frac{t_{r,\alpha}^2 s^2 v_{22}}{Y^2}$$

and  $t_{r,\alpha}$  is the  $\alpha$ -level deviate from the student's  $t$ -distribution based on  $r$  degrees of freedom.

*Proof.* Let  $Z = X - \theta Y$  and we have  $Z \perp\!\!\!\perp s^2$ . Thus,

$$\frac{Z/\sqrt{v_{11} - 2v_{12}\theta + v_{22}\theta^2}}{s} \sim t_r$$

where  $t_r$  is a student  $t$  variable with  $r$  degrees of freedom. Thus,

$$\mathbb{P}\left(\left(\frac{Z}{s\sqrt{v_{11} - 2v_{12}\theta + v_{22}\theta^2}}\right)^2 \leq t_{r,\alpha}^2\right) = 1 - \alpha$$

Solving for the quadratic equation

$$Z^2 = t_{r,\alpha}^2 s^2 (v_{11} - 2v_{12}\theta + v_{22}\theta^2)$$

It boils down to

$$(1 - g)\theta^2 - \left(\frac{2X}{Y} - \frac{2gv_{12}}{v_{22}}\right)\theta + X^2 - \frac{gv_{11}}{v_{22}} = 0$$

and the solution gives the upper and lower bound. □

**Example 6.1** (x-intercept estimation). Consider a simple linear model:

$$y - \bar{y} = \beta_0 + \beta(x - \bar{x})$$

To estimate the  $x$ -intercept, we let  $y = y_0$  and plug-in the usual LSE for  $\beta_0, \beta$ :

$$\hat{x}_0 = \frac{y_0 - \bar{y}}{\hat{\beta}} + \bar{x}$$

We apply Fieller's theorem to construct a confidence interval for the true  $x_0$ :

$$\begin{aligned} X &\leftarrow y_0 - \bar{y} \\ Y &\leftarrow \hat{\beta} \\ v_{11} &= \sigma^2 \left(1 + \frac{1}{n}\right) \\ v_{22} &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ v_{12} &= 0 \end{aligned}$$

Note  $v_{12} = 0$  follows from  $P_1 \perp P_x$ . Then Fieller's theorem gives

$$(x_L, x_U) = \bar{x} + \frac{1}{1-g} \left[ \frac{y_0 - \bar{y}}{\hat{\beta}} \mp \frac{t_{r,\alpha/2}}{\hat{\beta}} \sqrt{\sigma^2(1-g)\left(1 + \frac{1}{n}\right) + \frac{(y_0 - \bar{y})^2 \sigma^2}{\hat{\beta}^2 \sum_{i=1}^n (x_i - \bar{x})^2}} \right]$$

**Example 6.2** (Clinical trial). We are performing a prognostic factors interaction assessment in a 2-treatment trial with treatment  $T$  and  $C$ , along with a prognostic factor. The model is

$$\mathbb{E}y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

where  $x_2$  is a prognostic factor and

$$x_1 = \begin{cases} 1 & \text{if treatment } T \\ 0 & \text{if treatment } C \end{cases}$$

Thus,  $\mathbb{E}y_T = \beta_0 + \beta_1 + \beta_2 x_2 + \beta_3 x_2$ ,  $\mathbb{E}y_C = \beta_0 + \beta_2 x_2$  and

$$\mathbb{E}(y_T - y_C) = \beta_1 + \beta_3 x_2$$

We want to know  $x_2$  for which  $\mathbb{E}(y_T - y_C) > 0$ . That is, we want to estimate  $x_2 = -\frac{\beta_1}{\beta_3}$ . Let  $\hat{\beta}$  be the LSE and define

$$\tilde{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_3 \end{pmatrix}; \text{Cov}(\tilde{\beta}) = V = \begin{pmatrix} v_{11} & v_{13} \\ v_{31} & v_{33} \end{pmatrix}$$

Let  $w = \tilde{\beta}_1 + \theta \tilde{\beta}_3$ ,  $\theta = -\beta_1/\beta_3$ , then  $\mathbb{E}w = 0$  and  $\text{Var}(w) = (v_{11} + \theta^2 v_{33} + 2\theta v_{13})\sigma^2$ . We have

$$\frac{w}{\sqrt{\text{Var}(w)}} \sim \mathcal{N}(0, 1)$$

If  $\sigma^2$  is known, then solving for

$$(\hat{\beta}_1 + \theta \hat{\beta}_3)^2 = z_{\alpha/2}^2 (v_{11} + 2\theta v_{13} + \theta^2 v_{33})\sigma^2$$

gives the  $1 - \alpha$  confidence interval for  $\theta$ :

$$(\theta_L, \theta_U) = \frac{-w_2 \mp \sqrt{w_2^2 - 4w_1w_3}}{2w_1}$$

where

$$\begin{aligned} w_1 &= \hat{\beta}_1^2 - v_{11}z_{\alpha/2}^2 \\ w_2 &= 2\hat{\beta}_1\hat{\beta}_3 - 2v_{13}z_{\alpha/2}^2 \\ w_3 &= \hat{\beta}_3^2 - v_{33}z_{\alpha/2}^2 \end{aligned}$$

## 6.2 Lack of fitness test

Suppose we want to test

$$H_0 : \mathbb{E}Y = X\beta \text{ vs } H_a : \mathbb{E}Y \neq X\beta$$

with the following assumptions:

$$\begin{aligned} n_1 \text{ obs on } y \text{ at } x_1^T &= (x_{11}, \dots, x_{1k}) \\ n_2 \text{ obs on } y \text{ at } x_2^T &= (x_{21}, \dots, x_{2k}) \\ &\dots\dots\dots \\ n_g \text{ obs on } y \text{ at } x_g^T &= (x_{g1}, \dots, x_{gk}) \\ n &= \sum_{i=1}^g n_i \text{ is the total number of obs} \end{aligned}$$

In matrix form, we have

$$\mathbb{E}Y = WX\beta$$

where

$$W = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_2} & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{1}_{n_g} \end{pmatrix} X = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \dots \\ \mathbf{x}_g^T \end{pmatrix} = \begin{pmatrix} x_1 & x_2 & \dots & x_k \end{pmatrix}$$

Regressing  $y$  on  $x_1, \dots, x_k$  gives error sum of square **SSE** =  $y^T Q_{WX} y$ , pure error sum of

square **SSPE** =  $y^T Q_W y$  and lack of fitness sum of square **SSLoF** =  $y^T (Q_{WX} - Q_W) y$  where

$$Q_{WX} = I_n - WX (X^T W^T W X)^{-1} X^T W^T$$

$$Q_W = \begin{pmatrix} I_{n_1} - \frac{\mathbf{1}\mathbf{1}^T}{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & I_{n_2} - \frac{\mathbf{1}\mathbf{1}^T}{n_2} & \cdots & \mathbf{0} \\ \cdots & \cdots & \cdots & \cdots \\ \mathbf{0} & \mathbf{0} & \cdots & I_{n_g} - \frac{\mathbf{1}\mathbf{1}^T}{n_g} \end{pmatrix}$$

Then

$$\begin{aligned} \text{SSE} &= y^T (Q_{WX} - Q_W + Q_W) y \\ &= y^T (Q_{WX} - Q_W) y + y^T Q_W y \\ &= \text{SSLoF} + \text{SSPE} \end{aligned} \tag{6.2}$$

Since  $\mathcal{C}(WX) \subset \mathcal{C}(X)$ ,

$$\begin{aligned} (Q_{WX} - Q_W)Q_W &= Q_W - Q_W = 0 \\ (Q_{WX} - Q_W)(Q_{WX} - Q_W) &= Q_{WX} - Q_W \end{aligned}$$

This implies

$$y^T (Q_{WX} - Q_W) y \perp\!\!\!\perp y^T Q_W y$$

So **SSLoF**  $\perp\!\!\!\perp$  **SSPE** and

$$\begin{aligned} \text{SSLoF}/\sigma^2 &\sim \chi_{g-k}^2 \\ \text{SSPE}/\sigma^2 &\sim \chi_{n-g}^2 \end{aligned}$$

The test statistic for  $H_0 : \mathbb{E}Y = X\beta$  is

$$F = \frac{\text{SSLoF}/(g-k)}{\text{SSPE}/(n-g)} \sim F_{g-k, n-g} \tag{6.3}$$

### 6.3 Scheffe's method

Let  $A \in \mathbb{R}^{q \times p}$  and  $\text{rank}(A) = q$ . We want to find an interval estimate for  $\Phi = A\beta$ . For LSE, we have

$$\hat{\Phi} = A\hat{\beta} \sim \mathcal{N}_q(A\beta, \sigma^2 L)$$

where  $L = A(X^T X)^{-1} A^T$ . Thus,

$$\frac{1}{\sigma^2} (\hat{\Phi} - \Phi)^T L^{-1} (\hat{\Phi} - \Phi) \sim \chi_q^2$$

Besides,  $\frac{s^2(n-p)}{\sigma^2} = \frac{\hat{\sigma}^2(n-p)}{\sigma^2} \sim \chi_{n-p}^2$ . So,

$$\frac{(\hat{\Phi} - \Phi)^T L^{-1} (\hat{\Phi} - \Phi)}{qs^2} \sim F_{q, n-p}$$

and

$$1 - \alpha = \mathbb{P} \left( \frac{(\hat{\Phi} - \Phi)^T L^{-1} (\hat{\Phi} - \Phi)}{qs^2} \leq F_{q, n-p, \alpha} \right) \quad (6.4)$$

Solving the inequality in 6.4 for  $\Phi$  is non-trivial. However, Scheffe in 1953 provided a beautiful method to obtain confidence intervals for  $h^T \Phi$  where  $h$  is any  $q$ -dimensional vector.

**Theorem 6.2** (Scheffe's theorem). Let  $m = qs^2 F_{q, n-p, \alpha}$  and  $b = \hat{\Phi} - \Phi$ . The following holds:

$$\mathbb{P} (b^T L^{-1} b \leq m) = 1 - \alpha = \mathbb{P} \left( \forall h : h^T \Phi \in h^T \hat{\Phi} \pm \sqrt{m \cdot h^T L h} \right) \quad (6.5)$$

*Proof.*

$$\begin{aligned} 1 - \alpha &= \mathbb{P} (b^T L^{-1} b \leq m) \\ &=_{(*)} \mathbb{P} \left( \max_{h \neq 0} \frac{(h^T b)^2}{h^T L h} \leq m \right) \\ &= \mathbb{P} (\forall h, (h^T b)^2 \leq m \cdot h^T L h) \\ &= \mathbb{P} (\forall h, |h^T b| \leq \sqrt{m \cdot h^T L h}) \\ &= \mathbb{P} (\forall h : h^T \Phi \in h^T \hat{\Phi} \pm \sqrt{m \cdot h^T L h}) \end{aligned}$$

(\*) is due to the generalized Cauchy-Schwarz inequality 1.10. □

## 6.4 Tukey's $q$

This subsection includes Tukey's  $q$  distribution and multiple comparison tests for one-way ANOVA.

**Definition 6.1** (Tukey's  $q$ ). Let  $Z_1, \dots, Z_k$  and  $U$  be independent random variables with  $Z_i \sim \mathcal{N}(0, 1)$  and  $U \sim \chi_m^2(0)$ . Define

$$q = \max_{i \neq j} \frac{|Z_i - Z_j|}{\sqrt{U/m}}$$

Then  $q$  has a studentized range distribution with  $k$  and  $m$  degrees of freedom and write

$$q \sim q_{k, m}$$

**Example 6.3** (Tukey's pairwise comparison). In a one-way ANOVA, for  $j = 1, \dots, n$  and  $i =$

$1, \dots, k,$

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

$$\bar{y}_{i\cdot} = \mu + \alpha_i + \bar{\epsilon}_i.$$

$$\bar{y}_{j\cdot} = \mu + \alpha_j + \bar{\epsilon}_j.$$

We have (exercise)

$$\max_{i \neq j} \frac{\sqrt{n} |\bar{y}_{i\cdot} - \bar{y}_{j\cdot} - (\alpha_i - \alpha_j)|}{\hat{\sigma}} \sim q_{k, k(n-1)} \quad (6.6)$$

where  $\hat{\sigma}^2$  is the usual estimate of  $\sigma^2$ . It follows that

$$\begin{aligned} & \mathbb{P} \left( \alpha_i - \alpha_j \in \bar{y}_{i\cdot} - \bar{y}_{j\cdot} \pm \frac{\hat{\sigma}}{\sqrt{n}} q_{k, k(n-1), \alpha} \forall i \neq j \right) \\ &= \mathbb{P} \left( \frac{\sqrt{n} |\bar{y}_{i\cdot} - \bar{y}_{j\cdot} - (\alpha_i - \alpha_j)|}{\hat{\sigma}} \leq q_{k, k(n-1), \alpha} \forall i \neq j \right) \\ &= 1 - \alpha \end{aligned}$$

Next step is to construct sets of confidence intervals for **contrasts**

$$\sum_i c_i \alpha_i \quad \forall c \text{ s.t. } c^T \mathbf{1} = 0$$

in a one-way ANOVA. We have

$$\begin{aligned} & \mathbb{P} \left( \sum_i c_i \alpha_i \in \sum_i c_i \bar{y}_{i\cdot} \pm \frac{\hat{\sigma}}{\sqrt{n}} q_{k, k(n-1), \alpha} \left[ \sum_i |c_i|/2 \right], \forall c \right) \\ &= \mathbb{P} \left( \left| \sum_i c_i (\bar{y}_{i\cdot} - \alpha_i) \right| \leq \frac{\hat{\sigma}}{\sqrt{n}} q_{k, k(n-1), \alpha} \left[ \sum_i |c_i|/2 \right], \forall c \right) \\ &=_{(*)} \mathbb{P} \left( |\bar{y}_{i\cdot} - \alpha_i - (\bar{y}_{j\cdot} - \alpha_j)| \leq \frac{\hat{\sigma}}{\sqrt{n}} q_{k, k(n-1), \alpha}, \forall i, j \right) \\ &= 1 - \alpha \end{aligned}$$

where  $(*)$  comes from the following lemma.

**Lemma 14.** *Let  $a_1, a_2, \dots, a_k$  be real numbers. Then the following are equivalent*

•

$$|a_i - a_j| \leq b, \forall i, j$$

•

$$\left| \sum_{i=1}^k c_i a_i \right| \leq \frac{b}{2} \left( \sum_{i=1}^k |c_i| \right), \forall c \text{ s.t. } \sum_{i=1}^k c_i = 0$$



*Proof.* One direction is trivial: set  $c_i = 1$ ,  $c_j = -1$  and  $c_m = 0$  for  $m \neq i, j$ . The other direction is non-trivial but there is a simpler case in the exercise. □

## 6.5 Test of homogeneity

In ANOVA,  $k$  groups may have variances  $\sigma_i^2$ ,  $i = 1, \dots, k$ .

$$H_0 : \sigma_1 = \sigma_2 = \dots = \sigma_k = \sigma$$

$$H_1 : \sigma_i \neq \sigma_j, i \neq j$$

**Example 6.4** (Hartley's test). Let  $\hat{s}_i^2$  be the unbiased estimate of  $\sigma_i^2$ ,  $i = 1, \dots, k$ . The ratio

$$F_{\max} = \frac{\max_{1 \leq i \leq k} \hat{s}_i^2}{\min_{1 \leq i \leq k} \hat{s}_i^2} \quad (6.7)$$

is close to 1 if  $H_0$  holds, otherwise we reject  $H_0$ . This is Hartley's  $F_{\max}$  statistic.

**Example 6.5** (Levene's test). Given  $\{y_{ij} : i = 1, \dots, k; j = 1, \dots, n_i\}$ , we define

$$\begin{aligned} z_{ij} &= |y_{ij} - \bar{y}_{i\cdot}| \\ z_{i\cdot} &= \frac{1}{n_i} \sum_{j=1}^{n_i} z_{ij} \\ z_{..} &= \frac{1}{n} \sum_i^k \sum_j^{n_i} z_{ij} \end{aligned}$$

Then use

$$W = \frac{\sum_i^k n_i (z_{i\cdot} - z_{..})^2 / (k - 1)}{\sum_i \sum_j (z_{ij} - z_{i\cdot})^2 / (n - k)} \quad (6.8)$$

as the test statistic. This is called Levene's test.

**Example 6.6** (Brown-Forsythe test). If  $\bar{y}_{i\cdot}$  is replaced by  $\tilde{y}_{i\cdot}$ : the median in  $i^{th}$  group, then it is Brown-Forsythe test.

**Example 6.7** (Bartlett's test). In 1937, Bartlett proposed a modification of the likelihood ratio test (LRT) to test the homogeneity problem. Let  $s_i^2$  be the sample variance in  $i^{th}$  group,  $n = \sum_{i=1}^k n_i$  be the total number of observations and

$$s^2 = \frac{1}{n - k} \sum_{i=1}^k (n_i - 1) s_i^2$$

The Bartlett's test statistic is

$$\chi^2 = \frac{(n-k) \ln(s^2) - \sum_{i=1}^k (n_i - 1) \ln(s_i^2)}{1 + \frac{1}{3(k-1)} \left( \sum_{i=1}^k \left( \frac{1}{n_i - 1} \right) - \frac{1}{n-k} \right)} \quad (6.9)$$

The test statistic has approximately a  $\chi_{k-1}^2$  distribution.

## 6.6 Cook-Weisberg score test

This subsection is based on question 2 of the final exam of 250B. Suppose we would like to test whether the variance of a continuous response depends on the levels of several independent variables in a regression model, then what kind of statistical procedure should we use? I think the paper written by Cook and Weisberg [6] in 1983 provided an alternative and it deserves my effort of writing a subsection.

**Example 6.8** (Cook-Weisberg test, [6]). Assume the regression model is

$$Y = \beta_0 \mathbf{1} + X\beta + \epsilon$$

$$X_0 = (\mathbf{1}, X) \text{ is fixed.}$$

$$\epsilon \sim \mathcal{N}_n(0, \sigma^2 W)$$

where  $X \in \mathbb{R}^{n \times (p-1)}$ ,  $Y$  is the continuous response vector and

$$W = \begin{pmatrix} w_1 & 0 & \cdots & 0 \\ 0 & w_2 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & w_n \end{pmatrix}$$

and

$$w_i = w(\lambda, z_i)$$

where  $z_i$ 's are  $q_z$ -dimensional independent variables and  $\lambda$  is a  $q_\lambda$ -dimensional parameter to be tested. For example,  $w_i$  can be taken as

$$w_i = \exp\left(\sum_{j=1}^q \lambda_j z_{ij}\right)$$

$$w_i = \prod_{j=1}^q z_{ij}^{\lambda_j}$$

We want to test

$$H_0 : \lambda = 0 \text{ vs } H_1 : \lambda \neq 0 \quad (6.10)$$

WLOG, assume  $q_z = q_\lambda = q$  and

$$w(0, z_i) = 1, \quad W =_{H_0} I$$

Let  $\hat{\beta}_0$  and  $\hat{\beta}$  be the maximum likelihood estimation of  $\beta_0$  and  $\beta$  under  $H_0$ . Define

$$e_i = (y_i - \hat{\beta}_0 - \hat{\beta})$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n - p}$$

Next, define

$$U = \begin{pmatrix} e_1^2/\hat{\sigma}^2 \\ e_2^2/\hat{\sigma}^2 \\ \dots \\ e_n^2/\hat{\sigma}^2 \end{pmatrix} \quad w'_i = \begin{pmatrix} \frac{\partial w_i}{\partial \lambda_1} \\ \frac{\partial w_i}{\partial \lambda_2} \\ \dots \\ \frac{\partial w_i}{\partial \lambda_q} \end{pmatrix} \quad D = \begin{pmatrix} (w'_1)^T \\ (w'_2)^T \\ \dots \\ (w'_n)^T \end{pmatrix}$$

so  $D$  is a  $n \times q$  matrix. Finally, let

$$\bar{D} = D - \mathbf{1}\mathbf{1}^T D/n$$

be the  $n \times q$  matrix obtained from  $D$  subtracting its column means.

Then the **Cook-Weisberg score statistic** is defined as (exercise)

$$S = \frac{1}{2} U^T \bar{D} (\bar{D}^T \bar{D})^{-1} \bar{D}^T U = \frac{1}{2} U^T P_{\bar{D}} U \quad (6.11)$$

with the assumption that  $\bar{D}$  is of full rank. The asymptotic distribution of  $S$  is

$$S \rightarrow_d \chi_q^2$$

Thus, reject  $H_0$  if

$$S > \chi_{q,\alpha}^2$$

where  $\chi_{q,\alpha}^2$  is the upper  $\alpha$  quantile of a  $\chi_q^2$  random variable.

In the special case when  $q = 1$ , i.e.,  $z_i$  and  $\lambda$  are both scalar, an alternative to Cook-Weisberg score statistic is (exercise)

$$\tilde{S} = \frac{Y^T (I - V) A (I - V) Y}{Y^T (I - V) Y} \quad (6.12)$$

where  $V = P_{X_0} = X_0 (X_0^T X_0)^{-1} X_0^T$  and  $A = \text{diag}(w'_i)$ . By simultaneous diagonalization theorem 1.1, we have

$$S' = \frac{\sum_{i=1}^{n-p} \mu_i \chi_i^2}{\sum_{i=1}^{n-p} \chi_i^2} \quad (6.13)$$

where  $\chi_i^2$ 's are iid  $\chi_1^2$  variables and  $\mu_i$ 's are, at most,  $n - p$  nonzero eigenvalues of  $(I - V)A(I - V)$ . Denote the null distribution of  $\tilde{S}$  as  $\mathbb{P}_{\tilde{S}}$ , then reject null if

$$\tilde{S} > \mathbb{P}_{\tilde{S}, \alpha}$$

where  $\mathbb{P}_{\tilde{S}, \alpha}$  is the upper  $\alpha$  quantile of  $\mathbb{P}_{\tilde{S}}$ .

## 6.7 Exercises

1. (Final exam) Consider the standard linear model where  $\mathbb{E}y = X\beta$  and  $X$  is of full rank. Find a 95% confidence interval for a ratio of two given linear combinations of the model parameters  $\frac{c^T \beta}{d^T \beta}$  where  $c$  and  $d$  are given vectors.
2. Find a 95% confidence interval for the unknown change point in a two phase linear regression.
3. (Hoadley's formula, [1]) Consider a simple linear regression and let  $F = \hat{\beta}_1^2 \sum_i (x_i - \bar{x})^2 / s^2$ , the  $F$ -statistic for testing  $H_0 : \beta_1 = 0$  for a straight line. Using the notation of Section 6.1.5 in Seber and Lee [1], prove that

$$\tilde{x}_0 - \bar{x} = \frac{F}{F + (n - 2)} (\hat{x}_0 - \bar{x})$$

4. Prove the claim 6.6.
5. Suppose we have a standard one-way ANOVA with  $k$  groups and  $n$  observations per group with the usual notation.

(a) Show that the statistic

$$T_{k,n} = \max_{1 \leq i \leq j \leq k} \sqrt{n} |\bar{y}_{i\cdot} - \bar{y}_{j\cdot} - (\mu_i - \mu_j)| / \hat{\sigma}$$

is distributed as the studentized range distribution and identify its degrees of freedom.

(b) Find the asymptotic distribution of  $T_{k,n}$  as  $n$  goes to  $\infty$ .

6. (Studentized maximum modulus distribution, [7]) Let  $Z_1, \dots, Z_k$  be i.i.d.  $\mathcal{N}(0, 1)$  and  $U \sim \chi_m^2(0)$  and they are independent. Define

$$M = \max_{1 \leq i \leq k} \frac{|Z_i|}{\sqrt{U/m}}$$

We say that  $M$  has a studentized maximum modulus distribution and write

$$M \sim M_{k,m}$$

- (a) Find simultaneous confidence intervals for the set of all  $\mu_i = \theta + \alpha_i$  for an one-way ANOVA model with  $k$  groups and  $n$  observations in each group.

**Hint:** If  $Z = (Z_1, \dots, Z_k)$  has distribution

$$Z \sim \mathcal{N}_k(0, \Sigma)$$

Then we still have [8]

$$\mathbb{P} \left( \max_{1 \leq i \leq k} \frac{|Z_i|}{\sqrt{U/m}} \geq M_{k,m,\alpha} \right) \leq \alpha$$

where  $M_{k,m,\alpha}$  is the upper  $\alpha$  quantile of  $M_{k,m}$ .

- (b) Let  $a_1, a_2, \dots, a_k$  be a set of numbers. Is it true that  $\max_i |a_i| \leq c$  iff  $|\sum_{i=1}^k d_i a_i| \leq c \sum_{i=1}^k |d_i|$  for all numbers  $d_1, d_2, \dots, d_k$ ?
- (c) Use (b) and the studentized maximum modulus distribution to find a simultaneous set of confidence intervals for  $\sum_{i=1}^k d_i \mu_i$ .
7. Refer to the Cook and Weisberg's paper **Diagnostics for Heteroscedasticity in Regression** [6]. Derive the test statistic for testing heteroscedasticity using the score test for the situation described in the paper and see whether our results agree with equations (8), (9) and (10).
8. Let  $X \sim \chi_n^2$ . Find  $\mathbb{E}\sqrt{X}$  and hence find an unbiased estimate of  $\sigma$  in the standard linear model  $\mathbb{E}y = X\beta$  where  $X$  is  $n \times p$  with full column rank and error terms are independent, each with mean 0 and variance  $\sigma^2$ .
9. Suppose we have  $n$  data points with independent observations  $y_1, \dots, y_n$  from the simple linear model defined on a compact interval  $X$ :

$$y_i = \beta_0 + \beta_1 x_i + e_i, x \in X$$

$$e_i \sim \mathcal{N}(0, \sigma^2), i = 1, \dots, n$$

- (a) Find an unbiased estimator for  $\beta_0, \beta_1$  and  $\sigma$ .
- (b) Show that if  $n > 2$ , a  $100(1 - \alpha)\%$  confidence interval for  $\xi = \sigma/\beta_1$  is

$$\hat{\xi} \mp z_{\frac{\alpha}{2}} \hat{\xi} \sqrt{\frac{1}{xn - 4} + \frac{\hat{\xi}^2}{s_x^2}}$$

where  $s_x^2 = \sum_{i=1}^n (x_i - \bar{x})^2$ ,  $\hat{\xi} = \hat{\sigma}/\hat{\beta}_1$ , and  $z_{\frac{\alpha}{2}}$  is the upper  $\alpha$  percentile of the standard normal.

## 7 Shrinkage and Bayes Estimation

In this section, we study different alternatives to the ordinary least square estimation. Principal component regression, ridge regression, LASSO, James-Stein estimator and Bayes estimator are included.

### 7.1 Principal component regression

Consider the linear model

$$\mathbb{E}y = X\beta$$

and  $X$  is a  $n \times p$  matrix but the rank is  $r < p$ . In other words, we have a collinearity issue. Suppose the spectral decomposition of  $X^T X$  is

$$U^T X^T X U = D = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ 0 & 0 & \cdots & \lambda_p \end{pmatrix}$$

and  $\lambda_{r+1}, \dots, \lambda_p = 0$ . Re-write

$$\mathbb{E}y = X\beta = X U U^T \beta = X U_1 \gamma_1$$

where  $U_1$  is the first  $r$  columns of  $U$  and  $\gamma_1 = U_1^T \beta$ .  $Z = X U_1$  are called the principal components of  $X$ . So we regress  $y$  on  $z_1, \dots, z_r$  and note that

$$z_i^T z_j = \delta_{ij} u_i^T X^T X u_j$$

where  $u_i$  is the  $i^{th}$  column of  $U_1$  and  $z_i = X u_i$ . A screen plot is the plot of

$$\frac{\sum_{i=1}^j \lambda_i}{\sum_{i=1}^p \lambda_i}, j = 1, \dots, p$$

### 7.2 Ridge estimator

The mean square error (MSE) of  $\tilde{\beta}$  of  $\beta$  is

$$\text{MSE}(\tilde{\beta}) = \mathbb{E} \left( \tilde{\beta} - \beta \right) \left( \tilde{\beta} - \beta \right)^T = \text{Cov}(\tilde{\beta} + b b^T)$$

where  $b = \mathbb{E} \left( \tilde{\beta} - \beta \right)$  is called bias. Clearly, for LSE we have  $b(\hat{\beta}) = 0$  and  $\text{MSE}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$ .

**Definition 7.1** (Ridge estimator). The ridge estimator with parameter  $\kappa$  of  $\beta$  is defined as

$$b_\kappa = (X^T X + \kappa I)^{-1} X^T y, \quad \kappa > 0 \quad (7.1)$$

Let  $G_\kappa = (X^T X + \kappa I)^{-1}$ , then

$$\mathbb{E}b_\kappa = G_\kappa X^T X \beta = (I + \kappa(X^T X)^{-1})^{-1} \beta$$

Thus, by Woodbury identity,

$$\begin{aligned} \text{bias}(b_\kappa) &= -\kappa (I + \kappa(X^T X)^{-1})^{-1} (X^T X)^{-1} \beta \\ &= -\kappa (\kappa I + (X^T X))^{-1} \beta \\ &= -\kappa G_\kappa \beta \end{aligned}$$

and

$$\text{Var}(b_\kappa) = \sigma^2 G_\kappa X^T X G_\kappa$$

Thus,

$$\text{MSE}(b_\kappa) = G_\kappa (\sigma^2 X^T X + \kappa^2 \beta \beta^T) G_\kappa \quad (7.2)$$

It can be shown (exercise)

$$\text{Tr}(\text{MSE}(\hat{\beta}) - \text{MSE}(b_\kappa)) \geq 0 \quad (7.3)$$

for some  $\kappa$ . Let the SVD of  $X$  be

$$X_{n \times p} = U_{n \times p} D_{p \times p} V_{p \times p}^T$$

Let  $u_i$  be the  $i^{th}$  column of  $U$ , then

$$\begin{aligned} \hat{\theta} &= X \hat{\beta} \\ &= U U^T y \\ &= \sum_{i=1}^p u_i u_i^T y \\ \tilde{\theta} &= X b_\kappa \\ &= X (X^T X + \kappa I)^{-1} X^T y \\ &= \sum_{i=1}^p \frac{\lambda_i^2}{\lambda_i^2 + \kappa} u_i u_i^T y \end{aligned}$$

Since  $\frac{\lambda_i^2}{\lambda_i^2 + \kappa} < 1$ , so  $X b_\kappa$  shrinks  $X \hat{\beta}$ .

### 7.3 James-Stein estimator

For technical details, please see Chapter 11 of [9]. Suppose

$$y \sim \mathcal{N}_n(\mu, \sigma^2 I), \mu \in \mathcal{V}, \sigma^2 > 0$$

where  $\mathcal{V}$  is a  $p$ -dimensional vector space of  $\mathbb{R}^n$ . Only in this subsection, I use  $P$  to denote the projection operator  $\mathcal{P}_{\mathcal{V}}$  onto  $\mathcal{V}$  and  $Q$  is the abbreviation for  $Q_{\mathcal{V}}$ . The LSE is

$$\hat{\mu} = Py$$

and

$$\hat{\sigma}^2 = y^T Q y / (n - p)$$

However,  $\hat{\mu}$  over-estimates the square of the norm of  $\mu$ :

$$\|\hat{\mu}\|_2^2 = \|\mu\|_2^2 + \sigma^2 \text{Tr}(P) = \|\mu\|_2^2 + p\sigma^2 > \|\mu\|_2^2$$

Consider shrinking  $\hat{\mu}$  and consider estimators of the form

$$\hat{\tilde{\mu}} = \left(1 - c \frac{\hat{\sigma}^2}{\|\mu\|_2^2}\right) \hat{\mu}$$

where  $c$  is a constant to be determined. In 1960, James and Stein provided one choice of  $c$ :

**Definition 7.2** (James-Stein (JS) estimator). The JS estimator of  $\mu$  is

$$\hat{\tilde{\mu}} = \left(1 - \frac{(p-2)(n-p)}{n-p+2} \frac{\hat{\sigma}^2}{\|\mu\|_2^2}\right) \hat{\mu} \quad (7.4)$$

That is,

$$c = \frac{(p-2)(n-p)}{n-p+2}$$

They also provided the following theorem:

**Theorem 7.1** (Risk of JS estimator, [9]). The estimator  $\hat{\tilde{\mu}}$  is better than  $\hat{\mu}$ . Its risk function is

$$R(\hat{\tilde{\mu}}; (\mu, \sigma^2)) = p - \frac{(p-2)^2(n-p)}{n-p+2} \mathbb{E} \left( \frac{1}{p-2+2K} \right) < p = R(\hat{\mu}; (\mu, \sigma^2)) \quad (7.5)$$

where  $K$  has a Poisson distribution with mean  $\frac{\|\mu\|_2^2}{2\sigma^2}$ . In particular,

$$R(\hat{\tilde{\mu}}; (0, \sigma^2)) = \frac{2n}{n-p+2}$$

For proof of theorem 7.1, see [9].



## 7.4 Bayes estimator

As Dr. Wong mentioned, in 1970s and 1980s, there is a huge debate between Frequentism and Bayesianism. But starting at 1990s, it seems that Bayesianism beats Frequentism. Therefore, in this subsection we talk about Bayes estimation in linear models. Let  $f(y|\theta)$  be the conditional density of  $y$  and  $g(\theta)$  be the prior density of  $\theta$ . Then the likelihood function of  $\theta$  is

$$L(y|\theta) = \prod_{i=1}^n f(y_i|\theta)$$

The posterior distribution of  $\theta$  is

$$f(\theta|y) = \frac{f(y|\theta)g(\theta)}{\int_{\Theta} f(y|\theta)g(\theta)d\theta} \propto f(y|\theta)g(\theta)$$

Consider a linear model:

$$y = X\beta + \epsilon, \quad \epsilon \sim \mathcal{N}_n(0, \sigma^2 I)$$

In the following examples, we put different priors on the parameter  $\theta^T = (\beta^T, \log \sigma)$ . Recall

$$f(\theta|y) \propto f(y|\theta)g(\theta), \quad g(\theta) = g_1(\beta|\sigma)g_2(\sigma)$$

and

$$y|\theta \sim \mathcal{N}(X\beta, \sigma^2 I)$$

**Example 7.1** (Non-informative prior). We need formulas 7.6, 7.7 and 7.8 (exercises).

The non-informative prior is

$$g(\theta) = g_1(\beta)g_2(\sigma^2) \propto \frac{1}{\sigma}$$

Thus,

$$f(\theta|y) \propto \frac{1}{\sigma} \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left( -\frac{1}{2\sigma^2} \|y - X\beta\|_2^2 \right)$$

If interest is in  $\beta$  along (exercises),

$$\begin{aligned} f(\beta|y) &= \int_0^\infty \frac{1}{(\sqrt{2\pi})^n} \frac{1}{\sigma^{n+1}} \exp \left( -\frac{1}{2\sigma^2} \|y - X\beta\|_2^2 \right) d\sigma^2 \\ &= \frac{1}{2} \left( \frac{\|y - X\beta\|_2^2}{2} \right)^{-\frac{n}{2}} \Gamma\left(\frac{n}{2}\right) \\ &\propto \|y - X\beta\|_2^{-n} \\ &= \left( (n-p)s^2 + \|X\hat{\beta} - X\beta\|^2 \right)^{-\frac{n}{2}} \\ &\propto \left( 1 + \frac{(\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta)}{(n-p)s^2} \right)^{-\frac{n}{2}} \end{aligned}$$

In conclusion,

$$\beta|y \sim t_p((n-p), \hat{\beta}, (X^T X/s^2)^{-1})$$

**Definition 7.3** (Multivariate  $t$ -distribution). We call  $y$  has a Multivariate  $t$ -distribution and write

$$y \sim t_m(\nu, \mu, \Sigma)$$

if the density of  $y$  is

$$f_Y(y) = \frac{\Gamma(\frac{\nu+m}{2})}{\Gamma(\frac{\nu}{2})(\pi\nu)^{\frac{m}{2}} |\Sigma|^{1/2}} \left(1 + \frac{(y-\mu)^T \Sigma^{-1} (y-\mu)}{\nu}\right)^{-\frac{\nu+m}{2}}$$

For details, see p475 of [1].

**Example 7.2** (Conjugate prior). We need formula 7.9 (exercise).

We put a conjugate prior on  $\theta^T = (\beta^T, \sigma^2)$ , i.e.,

$$g(\beta, \sigma^2) = g(\beta|\sigma^2)g(\sigma^2)$$

where  $\beta|\sigma^2 \sim \mathcal{N}_p(m, \sigma^2 V)$  and  $\sigma^2$  has an inverted gamma density (7.10) with hyper-parameters  $(d, \frac{a}{2})$ . Thus,

$$f(\sigma^2) \propto \frac{1}{(\sigma^2)^{\frac{d+2}{2}}} \exp(-a/2\sigma^2)$$

The posterior of  $\theta$  is proportional to

$$\begin{aligned} f(\theta|y) &\propto f(y|\theta)g(\theta) \\ &\propto (\sigma^2)^{-\frac{n+d+p+2}{2}} \exp\left(-\frac{1}{2\sigma^2}(Q+a)\right) \end{aligned}$$

where

$$Q = (y - X\beta)^T (y - X\beta) + (\beta - m)^T V^{-1} (\beta - m)$$

If interest is in  $\beta$  alone, then (exercise)

$$\begin{aligned} f(\beta|y) &\propto \int_0^\infty (\sigma^2)^{-\frac{n+d+p+2}{2}} \exp\left(-\frac{1}{2\sigma^2}(Q+a)\right) d\sigma^2 \\ &\propto (Q+a)^{-\frac{d+n+p}{2}} \\ &\propto \left(1 + \frac{1}{n+d}(\beta - m_*)^T W_*^{-1} (\beta - m_*)\right)^{-\frac{d+n+p}{2}} \end{aligned}$$

For definition of  $m_*$  and  $W_*$ , see exercise 7.9. In conclusion,

$$\beta|y \sim t_p(n+d, m_*, W_*)$$

## 7.5 Exercises

1. Prove inequality 7.3 for some  $\kappa$ .
2. Let  $X$  be a  $p$ -dimensional random vector with covariance matrix  $\Sigma$ . Let  $u = (u_1, u_2, \dots, u_p)^T$  be a vector of principal components of  $X$ . Then  $u_i = a_i^T X$  for some vector  $a_i$  of length 1,  $i = 1, 2, \dots, p$ . Show that
  - (a)  $\text{Var}(a^T X) \leq \text{Var}(u_1)$ .
  - (b) If  $a^T X$  is uncorrelated with  $u_1, u_2, \dots, u_{i-1}$ , then  $\text{Var}(a^T X) \leq \text{Var}(u_1)$ .
3. Let  $\mu$  belong to  $\mathcal{V}$ , a  $p$ -dimensional vector space of  $\mathbb{R}^n$ , let  $y$  belong to  $\mathbb{R}^n$  and let  $P_{\mathcal{V}}$  be the orthogonal projection onto  $\mathcal{V}$ . If  $\hat{\mu} = P_{\mathcal{V}} y$  and  $\hat{\sigma}^2 = \|(I - P_{\mathcal{V}})y\|_2^2 / (n - p)$ , show that
  - (a)  $\mathbb{E}\|\hat{\mu}\|_2^2$  always over-estimates  $\|\mu\|_2^2$ .
  - (b)  $\mathbb{E}\|\hat{\mu} - \mu\|_2^2 / \sigma^2 = p$  and identify the distribution of  $\|\hat{\mu} - \mu\|_2^2$ .
  - (c)  $\mathbb{E}\hat{\sigma}^4 = (n - p + 2)\sigma^4 / (n - p)$ .
4. Recall that if  $X$  and  $Y$  are random variables with finite means and variances, then

$$\mathbb{E}(\mathbb{E}(X|Y)) = \mathbb{E}X$$

Use this result to show that if the conditional distribution of  $V|K \sim \chi_{p+2K}^2(0)$ , then

$$\mathbb{E}\left(\frac{1}{V}|K\right) = \frac{1}{p - 2 + 2K}$$

5. (Beta-Binomial) Let  $Y \sim \text{Bin}(n, \theta)$ , that is,  $\mathbb{P}(Y = y) = \binom{n}{y} \theta^y (1 - \theta)^{n-y}$ . Let  $\theta \sim \text{Beta}(\alpha, \beta)$ . Find the posterior of  $\theta$ .
6. (Normal sequence) Suppose the conditional density of  $y|\mu$  is  $\mathcal{N}(\mu, \sigma^2)$ . If  $\mu$  is univariate normal with mean  $\mu_0$  and variance  $\sigma_0^2$ , then
  - (a) Find the conditional density of  $\mu|y$ .
  - (b) Generalize the setting when we have a random sample  $y_1, \dots, y_n$ .
7. Prove the following:

- (a) Recall  $\Gamma(x) = \int_0^\infty \exp(-t)t^{x-1}dt$ , then show

$$\int_0^\infty \exp\left(-\frac{k}{x}\right)x^{-\nu-1}dx = k^{-\nu}\Gamma(\nu) \quad (7.6)$$

- (b) Let  $a$  and  $b$  be non-negative, then show

$$\int_0^\infty \exp\left(-\frac{a}{x^2}\right)x^{-b-1}dx = \frac{1}{2}a^{-b/2}\Gamma\left(\frac{b}{2}\right) \quad (7.7)$$

8. Show that

$$\|y - X\beta\|_2^2 = (n - p)s^2 + \|X\hat{\beta} - X\beta\|_2^2 \quad (7.8)$$

where  $s^2$  is usual estimate of  $\sigma^2$  and  $\hat{\beta}$  is the LSE.

9. Use binomial inversion theorem to show the equality:

$$\begin{aligned} (y - X\beta)^T(y - X\beta) + (\beta - m)^T V^{-1}(\beta - m) = \\ (\beta - m^*)^T V^{*-1}(\beta - m^*) + (y - Xm)^T (I + XVX^T)^{-1}(y - Xm) \end{aligned} \quad (7.9)$$

where  $V^* = (X^T X + V^{-1})^{-1}$  and  $m^* = V^*(X^T y + V^{-1}m)$ .

10. (Inverted gamma) Suppose  $X$  has an inverted gamma with density given by

$$f(x|a, b) = b^{-a} x^{-a-1} \exp -\frac{b}{x} / \Gamma(a) \quad (7.10)$$

where  $a > 0$  and  $b > 0$  and  $\Gamma$  is the the Gamma function.

- (a) Describe how values of  $a$  and  $b$  affect  $f(x|a, b)$  in terms of shape, skewness and symmetry.
  - (b) Find its mean and describe how this density is typically used in estimating paramters in a linear model using a Bayesian paradigm.
11. Suppose  $Z$  has a density proportional to  $z^{n/2-1} \exp -nz/2$  and conditional on  $Z$ ,  $X|Z = z$  is multivariate normal with mean 0 and covariance  $\frac{1}{z}I_k$ .
- (a) What is the density of  $X$ ?
  - (b) Determine the mean and covariance matrix of  $X$ .
12. (posterior mean of  $\beta$ , [1]) Using the non-informative prior for  $\theta^T = (\beta^T, \log \sigma)$ , show that the conditional posterior density  $f(\beta|y, \sigma)$  is multivariate normal. Hence deduce that the posterior mean of  $\beta$  is  $\hat{\beta}$ , the LSE.
13. Read the paper on **Bayes Estimation of Two-Phase Linear Regression Model** [10].
14. (Posterior mean of  $\sigma^2$ , [1]) Suppose that we use the non-informative prior for  $\theta$ .

- (a) Obtain an expression proportional to  $f(\beta, \sigma^2|y)$ .
- (b) Integrate out  $\beta$  to obtain

$$f(\sigma^2|y) \propto (\sigma^2)^{-\frac{\nu}{2}-1} \exp(-\frac{a}{\sigma^2})$$

where  $\nu = n - p$  and  $a = \|y - X\hat{\beta}\|_2^2/2$ .

- (c) Find the posterior mean of  $\sigma^2$ .

## 8 ANOVA Mixed Models

Random effect, fixed effect and mixed effect ANOVA models are included.

### 8.1 One-way ANOVA with random effects

Recall the one-way ANOVA with fixed effects:

$$\begin{aligned}y_{ij} &= \mu + \tau_i + \epsilon_{ij} \\ \epsilon_{ij} &\sim_{\text{iid}} \mathcal{N}(0, \sigma^2) \\ \tau_i &= \text{effect of treatment } i\end{aligned}$$

We want to test

$$\begin{aligned}H_0 &: \tau_1 = \tau_2 = \cdots = \tau_k \\ H_1 &: \text{They are not the same.}\end{aligned}$$

When interest is in more than  $k$  treatments, i.e., interest is in a population of treatment, then we sample  $k$  of them, suggesting that

$$\begin{aligned}\tau_i &\sim \mathcal{N}(0, \sigma_\tau^2) \\ \epsilon_{ij} &\sim \mathcal{N}(0, \sigma^2) \\ \tau_i &\perp\!\!\!\perp \epsilon_{ij} \quad \forall i \& j\end{aligned}$$

Here  $\tau_i$ 's are called random effects. Hypothesis of interest is

$$H_0 : \sigma_\tau^2 = 0 \quad H_1 : \sigma_\tau^2 > 0$$

and we have data  $\{y_{ij} : i = 1, \dots, k; j = 1, \dots, n\}$ .

$$\begin{aligned}\sum_i \sum_j (y_{ij} - \bar{y})^2 &= \sum_i \sum_j (y_{ij} - \bar{y}_{i\cdot} + \bar{y}_{i\cdot} - \bar{y})^2 \\ &= \sum_i \sum_j (y_{ij} - \bar{y}_{i\cdot})^2 + \sum_i \sum_j (\bar{y}_{i\cdot} - \bar{y})^2 \\ &= \text{SSE} + \text{SSTr}\end{aligned}$$

Since  $y_{ij} = \mu + \tau_i + \epsilon_{ij}$ , then

$$\begin{aligned} \text{SSTr} &= \sum_i^k \sum_j^n (\tau_i - \bar{\tau} + \bar{\epsilon}_{i\cdot} - \bar{\epsilon})^2 \\ \text{SSE} &= \sum_i^k \sum_j^n (\epsilon_{ij} - \bar{\epsilon}_{i\cdot})^2 \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}(\text{SSTr}) &= (k-1)(n\sigma_\tau^2 + \sigma^2) \\ \mathbb{E}(\text{SSE}) &= n(k-1)\sigma^2 \end{aligned}$$

**Example 8.1** (Point estimation and confidence interval). The point estimation for  $\hat{\sigma}_\tau^2$  is

$$\hat{\sigma}_\tau^2 = \frac{\text{MSTr} - \text{MSE}}{n}$$

where  $\text{MSTr} = \frac{\text{SSTr}}{k-1}$  and  $\text{MSE} = \frac{\text{SSE}}{n(k-1)}$ . Also,

$$\frac{\text{SSE}}{\sigma^2} \sim \chi_{k(n-1)}^2$$

then a  $(1 - \alpha)$  confidence interval for  $\sigma^2$  is

$$\frac{\text{SSE}}{\chi_{k(n-1), 1-\frac{\alpha}{2}}^2} < \sigma^2 < \frac{\text{SSE}}{\chi_{k(n-1), \frac{\alpha}{2}}^2} \quad (8.1)$$

and

$$\frac{(k-1)\text{MSTr}}{n\sigma_\tau^2 + \sigma^2} \sim \chi_{k-1}^2$$

indicates a  $(1 - \alpha)$  confidence interval for  $\sigma_\tau^2$

$$\frac{\text{SSTr}}{\chi_{k-1, 1-\frac{\alpha}{2}}^2} < \sigma_\tau^2 < \frac{\text{SSTr}}{\chi_{k-1, \frac{\alpha}{2}}^2} \quad (8.2)$$

Further,  $\text{MSE} \perp\!\!\!\perp \text{MSTr}$  implies

$$\frac{\text{MSTr}}{\text{MSE}} \frac{\sigma^2}{n\sigma_\tau^2 + \sigma^2} \sim F_{k-1, (n-1)k}$$

Thus, if we are interested in the ratio:

$$\theta = \frac{\sigma_\tau^2}{\sigma^2}$$

then a  $100(1 - \alpha)\%$  confidence interval is

$$\left( \frac{1}{F_{k-1, (n-1)k, 1-\frac{\alpha}{2}}} \frac{\text{MSTr}}{\text{MSE}} - 1 \right) \frac{1}{n} \leq \theta \leq \left( \frac{1}{F_{k-1, (n-1)k, \frac{\alpha}{2}}} \frac{\text{MSTr}}{\text{MSE}} - 1 \right) \frac{1}{n} \quad (8.3)$$

The intraclass correlation coefficient for effect  $\tau$  is defined as

$$\rho = \frac{\sigma_\tau^2}{\sigma_\tau^2 + \sigma^2}$$

So a  $100(1 - \alpha)\%$  confidence interval for  $\rho$  is

$$\frac{l}{1+l} \leq \rho \leq \frac{u}{1+u}$$

where  $l$  and  $u$  are confidence limits for  $\theta$  and noting that  $\rho = \frac{\theta}{\theta+1}$ .

## 8.2 Two factor ANOVA

Suppose we have two factors  $A$  and  $B$ , each with  $a$  and  $b$  levels. Then the ANOVA model is

$$y_{ijk} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + \epsilon_{ijk}$$

For fixed effects, we have, for example, constraint on  $\tau$ :

$$\sum_{i=1}^a \tau_i = 0$$

For random effects, we assume, say,

$$\tau_i \sim \mathcal{N}(0, \sigma_\tau^2), \quad i = 1, \dots, a$$

so  $\tau_i$ 's are independent but not identically distributed.

**Example 8.2** ( $A$  fixed,  $B$  fixed). see figure 2.

**Example 8.3** ( $A$  random,  $B$  random). see figure 3.

Case (i): A, B both fixed.

Factor	a	b	n	EMS
$\tau_i$	0	b	n	*
$\beta_j$	a	0	n	**
$(\tau\beta)_{ij}$	0	0	n	***
$\epsilon_{ijk}$	1	1	1	$\Delta$

$(*) E(MSA) = b n \frac{\sum \tau_i^2}{a-1} + \sigma^2$   
 $(**) E(MSB) = a n \frac{\sum \beta_j^2}{b-1} + \sigma^2$   
 $(***) E(MSAB) = n \frac{\sum \sum (\tau\beta)_{ij}^2}{(a-1)(b-1)} + \sigma^2$   
 $(\Delta) E(MSE) = \sigma^2$

Figure 2: Fixed effect.

Case (ii): A & B are random

Factor	a	b	n	EMS
$\tau_i$	1	b	n	$b n \sigma_\tau^2 + n \sigma_{\tau\beta}^2 + \sigma^2$
$\beta_j$	a	1	n	$a n \sigma_\beta^2 + n \sigma_{\tau\beta}^2 + \sigma^2$
$(\tau\beta)_{ij}$	1	1	n	$n \sigma_{\tau\beta}^2 + \sigma^2$
$\epsilon_{ijk}$	1	1	1	$\sigma^2$

Test  $G_\tau^2 = 0 \rightarrow \frac{MSA}{MSAB}$   
 $G_\beta^2 = 0 \rightarrow \frac{MSB}{MSAB}$   
 $G_{\tau\beta}^2 = 0 \rightarrow \frac{MSAB}{MSE}$

Figure 3: Random effect.

Case (iii): One R One F

Factor	a	b	n	EMS
$\tau_i$	0	b	n	$b n \frac{\sum \tau_i^2}{a-1} + n \sigma_{\tau\beta}^2 + \sigma^2$
$\beta_j$	a	1	n	$a n \sigma_\beta^2 + \sigma^2$
$(\tau\beta)_{ij}$	1	1	n	$n \sigma_{\tau\beta}^2 + \sigma^2$
$\epsilon_{ijk}$	1	1	1	$\sigma^2$

$H_{0A}: \tau_1 = \tau_2 = \dots = \tau_a: \frac{MSA}{MSAB}$   
 $H_{0B}: \sigma_\beta^2 = 0: \frac{MSB}{MSE}$   
 $H_{0C}: \sigma_{\tau\beta}^2 = 0: \frac{MSAB}{MSE}$

Figure 4: Mixed effect.

**Example 8.4** (A fixed, B random). see figure 4. The assumptions for  $(\tau\beta)$  are:

$$(\tau\beta)_{ij} \sim \mathcal{N}(0, \frac{a-1}{a} \sigma_{\tau\beta}^2)$$

$$\sum_{i=1}^a (\tau\beta)_{ij} = 0$$

### 8.3 Satterwaite approximation

The Satterwaite approximation (a.k.a. Welch-Satterwaite equation) is a powerful tool to derive approximate F-test in random effect models. It is also used to solve the famous Berhens-Fisher problem (exercise).

**Theorem 8.1** (Satterwaite approximation). Let  $U_1, U_2, \dots, U_k$  be independent  $\chi^2$  random variables with  $r_1, r_2, \dots, r_k$  degrees of freedom respectively. Define

$$U = \sum_{i=1}^k a_i U_i$$

where  $a_1, \dots, a_k$  is a sequence of positive numbers. Then we have

$$\frac{bU}{E(U)} \sim \chi_b^2$$

where

$$b = \frac{(EU)^2}{\sum_{i=1}^k \frac{(a_i EU_i)^2}{r_i}} = \frac{\left(\sum_{i=1}^k a_i r_i\right)^2}{\left(\sum_{i=1}^k a_i^2 r_i\right)}$$

Therefore, by method of moments, we can approximate  $b$  by

$$\hat{b} = \frac{(\sum_{i=1}^k a_i U_i)^2}{\sum_{i=1}^k (a_i^2 U_i^2 / r_i)} \quad (8.4)$$



*Proof.* Suppose we want to find constant  $v$  s.t.

$$\frac{vU}{\mathbb{E}(U)} \sim \chi_b^2$$

Take expectation on both sides:

$$\mathbb{E}\left(\frac{vU}{\mathbb{E}(U)}\right) = b \Rightarrow v = b$$

Take variance on both sides:

$$\begin{aligned} 2b &= \text{Var}(\chi_b^2) \\ &= \text{Var}\left(\frac{vU}{\mathbb{E}(U)}\right) \\ &= \frac{v^2}{(\mathbb{E}U)^2} \sum_{i=1}^k a_i^2 \text{Var}(U_i) \\ &= \frac{v^2}{(\mathbb{E}(U))^2} \sum_{i=1}^k (2r_i a_i^2) \end{aligned}$$

Replacing  $v$  with  $b$  gives the desired result. □

## 8.4 Exercises

1. (Behrens-Fisher problem) Show that the statistic for testing equality of means from two normal populations with different variances using sample size  $n_1$  and  $n_2$  has a  $t$ -distribution with degrees of freedom **approximately** equal to

$$\frac{(g_1 + g_2)^2}{g_1^2/(n_1 - 1) + g_2^2/(n_2 - 1)}$$

where  $g_i = s_i^2/n_i, i = 1, 2$  and  $s_1^2$  and  $s_2^2$  are the sample variances.

2. (Final exam) Suppose we have an **unbalanced** one-way ANOVA setting with the usual assumptions, except that we do not assume the variances of the responses from each group are equal. Derive a statistical procedure to test whether a given linear combination of the group means is zero.
3. (Three factor mixed ANOVA model) Suppose we conduct an experiment with 3 factors  $A, B$  and  $C$ . The two factors  $A, B$  are fixed and  $C$  is a random factor. The number of levels for each of the factors are  $a, b$  and  $c$ , respectively, with  $n$  observations per cell.

- (a) Find the expected mean square errors (EMS) for all the effects, see figure 2, 3 and 4.
- (b) Provide the test statistics for testing all effects in the model and state the rejection rules. This is known as approximate  $F$ -test.

## 9 Linear Mixed Models

Consider a linear mixed model (LMM):

$$Y = X\beta + Z\gamma + \epsilon$$

$$\mathbb{E}\epsilon = 0, \text{Cov}(\epsilon) = R$$

and

- $\beta$ : an unknown vector of fixed effects.
- $\gamma$ : a vector of random effects with mean 0 covariance  $D$ .
- $X, Z$  are known and  $\text{Cov}(\gamma, \epsilon) = 0$ .

We have

$$\mathbb{E}Y = X\beta; \text{Cov}(Y) = ZDZ^T + R := V$$

In this section, we discuss the Henderson's equation for solving MLE in LMMs. The Best Linear Unbiased Predictor (BLUP) theorem is also proved.

### 9.1 Henderson's equation

To find MLEs, we assume

$$\gamma \sim \mathcal{N}(0, D); \epsilon \sim \mathcal{N}(0, R)$$

Write down the likelihood function for  $Y$  and  $\gamma$  in two steps:

$$f_{y,\gamma}(y, \gamma) = f_{y|\gamma}(y) f_{\gamma}(\gamma)$$

$$\propto \frac{1}{|R|^{1/2}} \exp\left(-\frac{1}{2}(y - X\beta - Z\gamma)^T R^{-1}(y - X\beta - Z\gamma)\right) \times$$

$$\frac{1}{|D|^{1/2}} \exp\left(-\frac{1}{2}\gamma^T D^{-1}\gamma\right) \quad (9.1)$$

Differentiate  $\ln f_{y,\gamma}(y, \gamma)$  w.r.t.  $\beta$  and  $\gamma$  results (exercise)

$$\begin{pmatrix} X^T R^{-1} X & X^T R^{-1} Z \\ Z^T R^{-1} X & D^{-1} + Z^T R^{-1} Z \end{pmatrix} \begin{pmatrix} \beta \\ \gamma \end{pmatrix} = \begin{pmatrix} X^T R^{-1} Y \\ Z^T R^{-1} Y \end{pmatrix} \quad (9.2)$$

This is the set of Henderson's equations. It can be shown (exercise) the solution is

$$\hat{\beta}_w = (X^T V^{-1} X)^{-1} X^T V^{-1} y$$

$$\hat{\gamma} = D Z^T V^{-1} (y - X \hat{\beta}_w) \quad (9.3)$$

where

$$V^{-1} = (ZDZ^T + R)^{-1} = R^{-1}Z(D^{-1} + Z^TR^{-1}Z)^{-1}Z^TR^{-1}$$

and note that  $\hat{\beta}_w$  is unbiased for  $\beta$  and  $\mathbb{E}\hat{\gamma} = 0$ .

## 9.2 Best linear unbiased predictor (BLUP)

A recall of generalized linear model:

$$y = X\beta + \epsilon, \text{ rank}(X) = p$$

where  $\mathbb{E}(\epsilon) = 0$  and  $\text{Var}(\epsilon) = \sigma^2 V$ . Then we have the weighted LSE:

$$\hat{\beta}_w = (X^T V^{-1} X)^{-1} X^T V^{-1} y$$

**Definition 9.1** (Linear unbiased predictor (LUP)). Suppose  $u$  is a random variable with mean 0 and finite variance. A linear predictor  $d + a^T y$  of  $c^T \beta + u$  is **unbiased** iff

$$\mathbb{E}(d + a^T y) = \mathbb{E}(c^T \beta + u) = c^T \beta$$

or  $c^T \beta + u$  is **predictable** iff  $\exists$  an LUP of  $c^T \beta + u$ .

**Lemma 15.**  $c^T \beta + u$  is predictable iff

$$\exists a \text{ s.t. } c = X^T a$$

This implies  $d + a^T y$  is an LUP of  $c^T \beta + u$  iff

$$d = 0, c = X^T a$$

**Definition 9.2** (BLUP). Let  $\hat{u}$  be a predictor of  $u$ . If  $\hat{u}$  satisfies the following three requirements, we call it a best unbiased linear predictor (BLUP).

- $\hat{u}$  is a linear function of  $y$ .
- $\hat{u}$  is unbiased for  $u$ .

$$\mathbb{E}(\hat{u} - u) = 0$$

- If  $v$  is any LUP, then

$$\text{Var}(\hat{u} - u) \leq \text{Var}(v - u)$$

Recall the MSE of  $\hat{\theta}$  of  $\theta$  is

$$\text{MSE}(\hat{\theta}) = \mathbb{E}(\hat{\theta} - \theta)^2 = \text{Var}(\hat{\theta}) + (\text{Bias}(\hat{\theta}))^2$$

**Theorem 9.1** (BLUP). Consider a generalized linear model:

$$y = X\beta + \epsilon$$

A linear mixed model can be written as a generalized linear model where  $\epsilon$  is replaced by

$$Z\gamma + \epsilon$$

and  $V = ZDZ^T + R$ . Assume  $u$  is a variable with mean 0 and finite variance and

$$\mathbb{E}(\epsilon) = 0$$

$$\text{Cov}(\epsilon) = \sigma^2 V$$

$$\text{Cov}(\epsilon, u) = \sigma^2 K$$

where  $K$  is a known vector. Note we have no assumption of normality. Then  $c^T \hat{\beta}_w + \hat{u}$  has the smallest MSE among all LUP of  $c^T \beta + u$ , where

$$\begin{aligned}\hat{u} &= K^T V^{-1} (y - X \hat{\beta}_w) \\ \hat{\beta}_w &= (X^T V^{-1} X)^{-1} X^T V^{-1} y\end{aligned}$$

That is,  $c^T \hat{\beta}_w + \hat{u}$  is BLUP of  $c^T \beta + u$ .

*Proof.* Note that

$$c^T \hat{\beta}_w + \hat{u} = c^T \hat{\beta}_w + K^T V^{-1} (y - X \hat{\beta}_w)$$

has expectation

$$c^T \beta + 0 = \mathbb{E}(c^T \beta + u)$$

Further,

$$\begin{aligned}c^T \hat{\beta}_w + \hat{u} &= (c^T + K^T V^{-1} X) \hat{\beta}_w + K^T V^{-1} y \\ &= [(c^T - K^T V^{-1} X)(X^T V^{-1} X)^{-1} X^T V^{-1} + K^T V^{-1}] y \\ &= b^T y\end{aligned}$$

where  $b^T = (c^T - K^T V^{-1} X)(X^T V^{-1} X)^{-1} X^T V^{-1} + K^T V^{-1}$ . Thus,  $c^T \hat{\beta}_w + \hat{u}$  is **linear** in  $y$ . It is also **unbiased** for  $c^T \beta + u$ , we have

$$c = X^T b$$

If  $a^T y$  is any other LUP of  $c^T \beta + u$ , we have

$$c = X^T a$$

To derive  $\text{MSE}(a^T y)$ , note that

$$\text{Cov}((a - b)^T y, b^T y - u) = \sigma^2 (a - b)^T (Vb - K)$$

But

$$Vb - K = X(X^T V^{-1} X)^{-1} (c - X^T V^{-1} K)$$

Thus,

$$\begin{aligned} \text{Cov}((a - b)^T y, b^T y - u) &= \sigma^2 (a - b)^T (Vb - K) \\ &= \sigma^2 [(a - b)^T X] [(X^T V^{-1} X)^{-1} (c - X^T V^{-1} K)] \\ &= 0 \end{aligned}$$

Finally,

$$\begin{aligned} \text{MSE}(a^T y) &= \mathbb{E} (a^T y - (c^T \beta + u))^2 \\ &= \text{Var}(a^T y - u) \\ &= \text{Var}(a^T y - b^T y + b^T y - u) \\ &= \text{Var}(a^T y - u) + 2\text{Cov}((a - b)^T y, b^T y - u) + \text{Var}(b^T y - u) \\ &= \text{Var}(a^T y - u) + \text{Var}(b^T y - u) \\ &= \text{Var}(a^T y - u) + \text{Var}(b^T y - (c^T \beta + u)) \\ &= \text{Var}((a - b)^T y) + \text{MSE}(b^T y) \\ &\geq \text{MSE}(b^T y) \end{aligned} \tag{9.4}$$

□

Recall the solution to Henderson's equation 9.3, we have

$$\hat{\gamma} = DZ^T V^{-1} (y - X\hat{\beta})$$

and note

$$\begin{pmatrix} \gamma \\ y \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ X\beta \end{pmatrix}, \begin{pmatrix} D & ? \\ ? & V \end{pmatrix} \right)$$

where

$$? = \text{Cov}(\gamma, y) = DZ^T$$

Thus (exercise),

$$\gamma|y \sim \mathcal{N}_q (DZ^T V^{-1} (y - X\beta), D - DZ^T V^{-1} ZD)$$

The **BLUP** of  $\hat{\gamma}$  (elementwise) is indeed

$$\widehat{\mathbb{E}(\gamma|y)} = DZ^T V^{-1}(y - X\hat{\beta}_w) \quad (9.5)$$

**Example 9.1** (BLUP for one-way ANOVA). Assume for  $i = 1, 2, 3$  and  $j = 1, 2$ ,

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

$$y = \mathbf{1}_6 \mu + Z\alpha + \epsilon$$

$$R = \text{Cov}(\epsilon) = \sigma^2 I_6$$

$$D = \text{Cov}(\alpha) = \sigma_\alpha^2 I_3$$

$$V = ZDZ^T + R$$

where  $\alpha = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}$ ,  $\mathbf{1}_6$  is a 6-dimensional vector of all ones and

$$Z = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$$

In this case, the BLUP for  $\alpha$  is

$$\hat{\alpha} = DZ^T V^{-1}(y - X\hat{\mu}_w)$$

where

$$\hat{\mu}_w = \frac{\mathbf{1}^T V^{-1} y}{\mathbf{1}^T V^{-1} \mathbf{1}}$$

### 9.3 Restricted maximum likelihood (REML)

Let the linear model be

$$Y = X\beta + \epsilon$$

$$\mathbb{E}\epsilon = 0, \text{Cov}(\epsilon) = V(\theta)$$

and

- $\beta$ : an unknown vector of fixed effects.
- $V(\theta)$ : the variance component depending on a parameter  $\theta$ .

- $\theta$ : an unknown vector of parameters.

Suppose we are only interested in estimating  $V(\theta)$ . Thus,  $\beta$  is the **nuisance parameter** and will cause extra variance when estimating  $V := V(\theta)$ . In 1974, David Harville proposed the so-called restricted maximum likelihood (REML) method [11]. In the original paper, David omitted details obtaining the REML. However, Lynn LaMotte provided a detailed derivation in 2007 [12]. In this subsection, I use LaMotte's derivation to get REML.

Let  $A$  be a matrix such that  $A^T X = 0$ , a particular choice of  $A$  would be

$$A = Q_X = I - X(X^T X)^{-1} X^T$$

Then we have

$$\mathbb{E}(AY) = 0, \text{Var}(AY) = AVA^T$$

and

$$L(\theta|y) = (2\pi)^{-\frac{n-p}{2}} |A^T V A|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} y^T A (A^T V A)^{-1} A^T y\right\} \quad (9.6)$$

This is the restricted maximum likelihood function.

**Theorem 9.2** (Re-parametrization of REML, [11], [12]). The REML of 9.6 can be re-written as

$$L(\theta|y) = (\text{Const.}) |V|^{-\frac{1}{2}} |X^T V^{-1} X|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} (y - X\hat{\beta}_w)^T V^{-1} (y - X\hat{\beta}_w)\right\} \quad (9.7)$$

with  $\hat{\beta}_w = (X^T V^{-1} X)^{-1} X^T V^{-1} y$  being the weighted LSE. In the special case  $A^T A = I$ ,

$$L(\theta|y) = (2\pi)^{-\frac{n-p}{2}} |X^T X|^{\frac{1}{2}} |V|^{-\frac{1}{2}} |X^T V^{-1} X|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} (y - X\hat{\beta}_w)^T V^{-1} (y - X\hat{\beta}_w)\right\}$$

LaMotte used the following two lemmas to show 9.7. Here I provide an easier proof of the first lemma and copy the original proof of the second lemma, which is first shown in Searle, 1979 [12].

**Lemma 16.** *If  $V$  is an  $n \times n$  positive definite matrix and  $X$  and  $A = Q_X$ , then*

$$V^{-1} = V^{-1} X (X^T V^{-1} X)^{-1} X^T V^{-1} + A (A^T V A)^{-1} A^T$$

*Proof.* Left-multiplying  $V^{\frac{1}{2}}$  and then right-multiplying  $V^{\frac{1}{2}}$  on both sides, it is enough to show

$$I = P_{V^{-\frac{1}{2}} X} + P_{V^{\frac{1}{2}} A}$$

But this follows directly from the fact  $V = V(\theta)$  is non-singular and

$$A^T X = 0, \mathcal{C}(X) \cup \mathcal{C}(A) = \mathbb{R}^n$$

□

**Lemma 17.** *If  $V$  is an  $n \times n$  positive definite matrix, and  $(X, A)$  is an  $n \times n$  matrix with full column rank, and  $A^T X = 0$ , then*

$$|V| = \frac{|A^T V A| |X^T X|}{|A^T A| |X^T V^{-1} X|}$$

*Proof.* Since  $(X, A)$  is a square matrix,

$$\begin{aligned} |A^T A| |X^T X| |V| &= \left| \begin{pmatrix} X^T X & 0 \\ 0 & A^T A \end{pmatrix} \right| |V| \\ &= |(X, A)^T V (X, A)| \\ &= \left| \begin{pmatrix} X^T V X & X^T V A \\ A^T V X & A^T V A \end{pmatrix} \right| \\ &=_{(*)} |A^T V A| |X^T V X - X^T V A (A^T V A)^{-1} A^T V X| \\ &=_{(@)} |A^T V A| |X^T (X (X^T V X)^{-1} X^T) X| \\ &= |A^T V A| \frac{|X^T X|^2}{|X^T V^{-1} X|} \end{aligned}$$

where  $(*)$  follows from 1.4 and  $(@)$  is due to 16. □

Combine these two lemmas we can get theorem 9.2 easily (exercise).

## 9.4 Exercises

1. Find  $\mathbb{E}(\gamma|y)$  under normality assumption.
2. Derive Henderson's equation 9.2.
3. Use binomial inversion theorem to derive 9.3.
4. Prove 9.7.



## 10 Spline Regression

In practice, we encounter situations where the relation between the co-variate  $x$  and the response  $y$  is non-linear and we do not have an explicit parametric form of such non-linear relation. Then we have to model the association between  $y$  and  $x$  using nonparametric regression. For example, Storey et al. [13] applied basis regression; Trapnell et al. [14] considered the GAM with the Tobit likelihood; Ren and Kuan [15] applied the GAM with Bayesian shrinkage dispersion estimates; Van den Berge et al. [16] proposed tradeSeq using the spline-based GAM. More recently, Song and Li [17] proposed the PseudotimeDE method, which fixes the p-value calibration issue in tradeSeq and also uses the spline-based GAM with spline functions. Additionally, Bacher et al. [18] used a piecewise linear model, which is more restrictive than the GAM.

In this section, we introduce the concept of spline regression.

### 10.1 Regression splines

**Definition 10.1** (Spline function). Given the data  $(x_i, y_i), i = 1, \dots, N$  and the breakpoints  $b_0, \dots, b_K$ , the spline function  $s(x)$  of order  $M$  is a piece-wise polynomial function satisfying the following conditions:

- For  $x \in [b_i, b_{i+1}]$ ,  $s(x)$  is a polynomial function of order (at most)  $M$ .
- The function has continuous derivatives up to order  $M-2$  at breakpoints, i.e.,  $\frac{\partial^m}{\partial x^m} s(x)|_{x=b_i} = \frac{\partial^m}{\partial x^m} s(x)|_{x=b_{i+1}}$  for  $i = 1, \dots, K-2$  and  $m = 0, 1, 2, \dots, M-2$ .
- The breakpoints  $b_0$  equals to the lower bound of the domain of  $x$  and  $b_K$  equals to the upper bound of the domain of  $x$ .

We usually refer to the splines described above as regression splines.

There are some particular cases [19]:

- If  $M = 4$ , then  $s(x)$  is called a cubic spline.
- If  $M = 2$ , then  $s(x)$  is called a linear spline.
- If we restrict  $s(x)$  such that  $\frac{\partial^2}{\partial x^2} s(x)|_{x=b_0} = \frac{\partial^2}{\partial x^2} s(x)|_{x=b_K} = 0$ , then  $s(x)$  is called a natural spline.
- If we let 0 in the above to be any other values specified by users, then  $s(x)$  is called a clamped spline.

- If  $s(b_0) = s(b_K)$ ,  $s'(b_0) = s'(b_K)$  and  $s''(b_0) = s''(b_K)$  then  $s(x)$  is referred as a periodic spline.

For cubic splines, we have  $K$  regions for regression, resulting in a total of  $4 \times K$  parameters; however, at each knot (breakpoint), we have 3 constraints. Hence, the degrees of freedom is  $4K - 3(K - 1) = K + 3$ . As pointed out in [1], a restricted least squares approach using the constraints at the knots is cumbersome, and more parameters need to be estimated than the "minimum"  $K + 4$ . Fortunately, any cubic spline with knots  $b_i$  has a unique representation in the form

$$s(x) = \sum_{h=0}^3 \alpha_h x^h + \sum_{j=1}^{K-1} \beta_j (x - b_j)_+^3 \quad (10.1)$$

where  $u_+ = \max(u, 0)$ . The above is not recommended for computational use because if  $x$  is much larger than  $b_j$  then the buildup of powers of large numbers leads to the ill-conditioning matrices.

A quote from [20]: "When we fit a spline, where should we place the knots? The regression spline is most flexible in regions that contain a lot of knots, because in those regions the polynomial coefficients can change rapidly. Hence, one option is to place more knots in places where we feel the function might vary most rapidly, and to place fewer knots where it seems more stable. While this option can work well, in practice it is common to place knots in a uniform fashion." There has been extensive literature on treating position of knots as free variable, and here are the comments from Dr. Telesca (re-organized, personal communication):

"The selection of knots position in spline smoothing usually finds applications in the recovery of functions that have different degrees of smoothness in a specific evaluation domain. Equivalently, one could have data with different noise levels in different intervals of the evaluation domain.

Some classical applications have been discussed by:

1) <https://www.stat.cmu.edu/~kass/papers/bars.pdf>

More recently, in the Gaussian Process literature:

2) <https://bobby.gramacy.com/prepo/gra2006-01.pdf>

There is a large body of literature on free-knots splines. Some recent applications to functional data analysis have been discussed by Gervini (<https://cpb-us-w2.wpmucdn.com/sites.uwm.edu/dist/6/111/files/2016/04/free-knot-sm-1atjte.pdf>).

Finally, in the literature on advanced computation, free-knots basis functions have been used to evaluate elementary functions in a numerically efficient way.

3) <https://ieeexplore.ieee.org/document/7926984>

This final paper, poses an interesting optimization problem, as knots for the radial basis need to be positioned optimally, and the value of basis coefficients are constrained to be powers of 2 for efficient fixed-point implementations."

Further, people are adopting the frequentist framework to optimize the locations of knots, see e.g. [21–24]. For [24], the authors novelly apply the particle swarm optimization (PSO) to solve for a high dimensional non-convex knot-positioning problem.

## 10.2 B-splines

This subsection is mainly based on the paper by [24], the appendix of chapter 5 in [25], section 7.2 in [1] and the 1996 monograph [26]. An  $M^{\text{th}}$ -order basis spline is an  $(M - 1)^{\text{th}}$ -degree piecewise polynomial which is positive in the interior of a domain of  $M$  intervals spanned by  $M + 1$  consecutive knots, and ZERO elsewhere. To define the complete basis, one needs to add  $2(M - 1)$  more knots at the lower and the upper bound of the domain. Denote these knots by

$$b_{-(M-1)}, b_{-(M-2)}, \dots, b_{-1}, b_0, b_1, \dots, b_K, b_{K+1}, b_{K+2}, \dots, b_{K+M}$$

where  $b_0 = a$  and  $b_{K+1} = b$  are the lower and upper bounds, they must satisfy

$$b_{-(M-1)} \leq b_{-(M-2)} \leq \dots \leq b_0 = a \text{ and } b = b_{K+1} \leq b_{K+2} \leq \dots \leq b_{K+M}$$

In literature, it is common to relabel the knots as  $\tau_{j+M} = b_j$  for all  $j$ , so that we have knots  $\tau_1, \dots, \tau_{K+2M}$ . Then B-splines are defined in terms of  $x, \tau, K$  and  $M$ .

**Definition 10.2** (B-spline). For  $m \leq M$ , the family of B-splines are defined as

$$B_{j,m}(x) = (\tau_{j+m} - \tau_j) \sum_{h=j}^{j+m} \frac{(x - \tau_h)_+^{m-1}}{\prod_{s=j, s \neq h}^{j+m} (\tau_h - \tau_s)} \quad (10.2)$$

for  $j = 1, 2, \dots, K + 2M$  where  $B_{j,m}(x)$ , known as  $j^{\text{th}}$  basis function of order  $m$ , is positive in the interval  $(\tau_j, \tau_{j+m})$  and has a single local maximum.

In 1972, de Boor provided a recursion formula [27, 28] for calculating  $B_{j,m}(x)$ ,

$$B_{j,1} = \mathbb{I}(\tau_j \leq x < \tau_{j+1})$$

for  $j = 1, 2, \dots, K + 2M - 1$  and

$$B_{j,m}(x) = \omega_{j,m}(x) B_{j,m-1}(x) + \gamma_{j+1,m}(x) B_{j+1,m-1}(x) \quad (10.3)$$

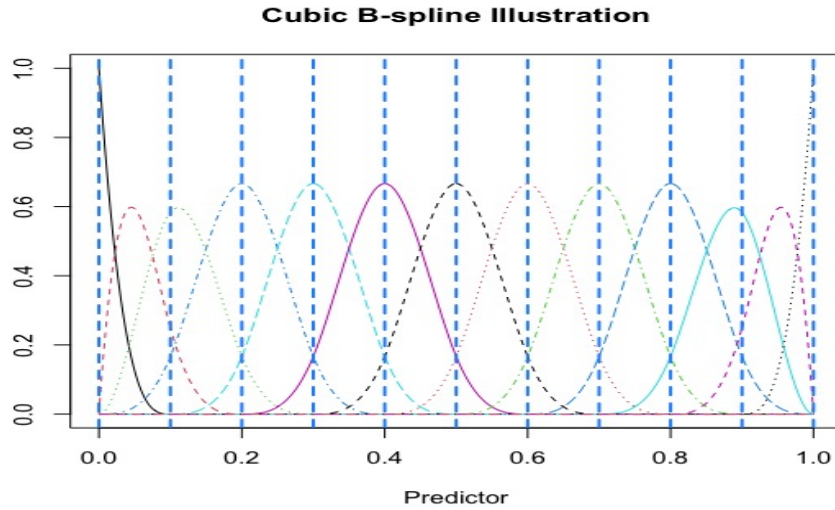
where

$$\omega_{j,m}(x) = \begin{cases} \frac{x - \tau_j}{\tau_{j+m-1} - \tau_j} & \text{if } \tau_{j+m-1} \neq \tau_j \\ 0 & \text{if } \tau_{j+m-1} = \tau_j \end{cases}, \quad \gamma_{j,m}(x) = \begin{cases} 1 - \omega_{j,m}(x) & \text{if } \tau_{j+m-1} \neq \tau_j \\ 0 & \text{if } \tau_{j+m-1} = \tau_j \end{cases}$$

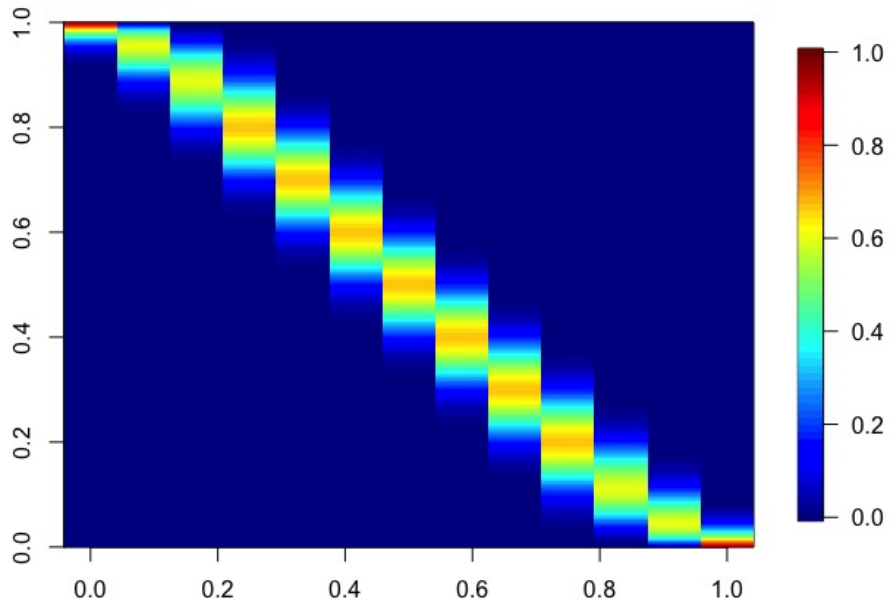
for  $j = 1, 2, \dots, K + 2M - m$ . In the special case that  $M = m$ , we have a piecewise polynomial (or order M) representation of  $s(x)$ , i.e.,

$$s(x) = \sum_{j=1}^{K+M} \alpha_j B_{j,M}(x) \quad (10.4)$$

where  $\alpha_j$ 's are regression coefficients. Figure 6 illustrates the B-splines for  $K = 4$  (cubic B-splines) with 11 breakpoints  $0, 0.1, 0.2, \dots, 1.0$ . We can see that every basis spline has positive values on exactly  $K = 4$  intervals. The fact that each basis spline is locally supported suggests that the design matrix  $\mathbf{X}$  has a band-like structure and  $\mathbf{X}$  has dimension  $N$ .



**Figure 5:** Cubic B-spline.



**Figure 6:** Design matrix of Cubic B-spline.

We have the following properties of B-splines.

**Theorem 10.1** (Properties of B-splines [25, 26]). Suppose  $B_{j,M}(x)$  is an order- $M$  B-spline defined in 10.2 using de Boor's recursion formula, then the following properties hold.

1. (Local support property I) For  $x \notin [\tau_j, \tau_{j+M}]$ , we have  $B_{j,M}(x) = 0$ .
2. (Local support property II) For  $x \in (\tau_j, \tau_{j+M})$ , we have  $B_{j,M}(x) > 0$ .
3. (Partition of unity) For any  $x \in (b_0, b_{K+1}) = (a, b)$ , we have  $\sum_{j=1}^{K+M} B_{j,M}(x) = 1$ .
4. The basis spline  $B_{j,M}(x)$  is a piece-wise polynomial of order  $M$  (degree  $M-1$ ) on  $[b_0, b_{K+1}] = [a, b]$ , with breaks only at the knots  $b_1, b_2, \dots, b_K$ .
5. An order- $M$  B-spline function can be represented as the density function of a convolution of  $M$  uniform random variables.
6. At a knot  $b_j$  of multiplicity  $r_j$ , the basis function  $B_{j,M}$  is  $\mathbb{C}^{M-r_j-1}$  continuous. Therefore, increasing multiplicity decreases the level of continuity, and increasing degree increases continuity.
7. Except for the case  $m = 1$ ,  $B_{j,m}(x)$  attains exactly one maximum value.
8. (Derivatives) The derivative of B-spline function is

$$\frac{dB_{j,m}(x)}{dx} = \frac{m-1}{\tau_{j+m-1} - \tau_j} B_{j,m-1}(x) - \frac{m-1}{\tau_{j+m} - \tau_{j+1}} B_{j+1,m-1}(x) \quad (10.5)$$

Let  $B_{j,m}^{(p)}(x)$  be the  $p^{th}$  derivative. Then repeated differentiation of the above produces the general formula

$$B_{j,m}^{(p)} = (m-1) \left( \frac{B_{j,m-1}^{(p-1)}(x)}{\tau_{j+m-1} - \tau_j} - \frac{B_{j+1,m-1}^{(p-1)}(x)}{\tau_{j+m} - \tau_{j+1}} \right) \quad (10.6)$$

9. (Butter's formula) Similar to de Boor's formula, we have a counterpart for calculating the derivative given by [29]

$$B_{j,m}^{(p)} = \frac{m-1}{m-1-p} \left( \omega_{j,m}(x) B_{j,m-1}^{(p)}(x) + \gamma_{j+1,m}(x) B_{j+1,m-1}^{(p)}(x) \right) \quad (10.7)$$

where  $\omega$  and  $\gamma$  are defined as before.

10. (Relation to Bernstein polynomials) A knot vector of the form ( $a = 0, b = 1$ )

$$\tau = \{\underbrace{0, \dots, 0}_m, \underbrace{1, \dots, 1}_m\}$$

generates the Bernstein polynomials of degree  $m-1$  (order  $m$ ).

*Proof.* Part 1:

$$\begin{aligned}
x \notin [\tau_j, \tau_{j+M}] &\Rightarrow B_{j,1} = 0, B_{j+1,1} = 0, \dots, B_{j+M-1,1} = 0 \\
&\Rightarrow B_{j,2} = 0, B_{j+1,2} = 0, \dots, B_{j+M-2,2} = 0 \\
&\dots \\
&\Rightarrow B_{j,M-1} = 0, B_{j+M-(M-1),M-1} = 0 \\
&\Rightarrow B_{j,M} = 0
\end{aligned}$$

Part 2: change  $\notin$  to  $\in$  and then all equal signs become " $>$ ".

Part 3: If  $M = 1$ , we trivially have

$$\sum_{i=1}^{K+1} B_{i,1}(x) = 1$$

If  $M = 2$ , we have

$$\begin{aligned}
\sum_{i=1}^{K+2} B_{i,2}(x) &= \sum_{i=1}^{K+2} (\omega_{i,2}(x) B_{i,1}(x) + \gamma_{i+1,1}(x) B_{i+1,1}(x)) \\
&= \sum_{i=2}^{K+1} (\omega_{i,2}(x) + \gamma_{i,2}(x)) B_{i,1}(x) \\
&= \sum_{i=2}^{K+1} B_{i,1}(x) \\
&= 1
\end{aligned}$$

where the second equality comes from the fact  $\omega_{1,2}(x) = 0, \gamma_{K+3,2}(x) = 0$ .

Next, for the general  $M$ , we have

$$\begin{aligned}
\sum_{i=1}^{K+M} B_{i,2}(x) &= \sum_{i=1}^{K+M} (\omega_{i,m}(x) B_{i,m-1}(x) + \gamma_{i+1,m}(x) B_{i+1,m-1}(x)) \\
&= \sum_{i=2}^{K+M-1} B_{i,m-1}(x) + 0 + 0 \\
&\dots \\
&= \sum_{i=M}^{K+1} B_{i,1}(x) \\
&= 1
\end{aligned}$$

Part 4: TBD. Part 5: TBD.

Part 6: TBD (actually, I don't know how to prove this and I cannot find a rigorous proof online).

Part 7: TBD. Part 8: TBD. Part 9: TBD. Part 10: TBD. □

Finally, we have the following analytical properties of B-splines (from Dr. Telesca's 285 notes).

**Theorem 10.2** (Spline Approximations). Let  $g \in \mathbb{C}^\alpha[0, 1]$ , be the set of functions on  $[0, 1]$  with  $\alpha$  continuous derivatives. Consider a set of  $J$  spline basis functions of order  $q \geq \alpha > 0$ , s.t.

$$g(x) \approx \sum_{j=1}^J \theta_j B_j(x) = \boldsymbol{\theta}^T B_J$$

where  $B_J$  is the vector of B-splines. There exists a constant  $M(q, \alpha) < \infty$  that depends only on  $\alpha$  and  $q$ , s.t.  $\forall g \in \mathbb{C}^\alpha[0, 1]$ , we can find  $\boldsymbol{\theta} \in \mathbb{R}^J$ , with  $\|\boldsymbol{\theta}\|_\infty < \|g\|_{\mathbb{C}^\alpha}$ , s.t.

$$\|g - \boldsymbol{\theta}^T B_J\|_\infty \leq M(\alpha, q) J^{-\alpha} \|g\|_{\mathbb{C}^\alpha} \quad (10.8)$$

That is, the uniform distance between a function  $g$  and the closest spline approximation of order  $q \geq \alpha$  is of the order  $J^{-\alpha}$ .

### 10.3 Smoothing splines

Let  $W_2[a, b]$  be the space of all smooth functions defined on  $[a, b]$ . Here "smooth" means if  $g \in W_2[a, b]$ , then  $g$  is second-times differentiable (almost everywhere), with the second-derivative being Lebesgue integrable over  $[a, b]$ . Given a set of points  $\{(x_i, y_i)\}, i = 1, 2, \dots, n$  and  $a = \min_i(x_i)$  and  $b = \max_i(x_i)$ , we have the following results.

**Theorem 10.3** (NCS [30], [31]). The solution to the following optimization problem is a natural cubic spline (NCS) with knots at  $n$  points  $x_1, x_2, \dots, x_n$ .

$$\min_{f \in W_2[a, b]} \left( \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda J(f) \right) \quad (10.9)$$

where  $\lambda$  is a smoothing parameter and  $J(f) = \int_a^b f''(x)^2 dx$ .

We write  $\text{RSS}(f, \lambda)$  for the objective function in theorem 10.3 and it is called the penalized residual sum of squares.

*Proof.* Let  $g$  be a NCS with knots  $x_1, \dots, x_n$  and  $f \in W_2([a, b])$  be an interpolant of  $\{x_i, y_i\}$  other than  $g$ . Define  $h(x) = f(x) - g(x)$  and we seek an expression for  $J(f)$  using  $J(g)$ .

$$\begin{aligned} \int_{x_1}^{x_n} f''(x)^2 dx &= \int_{x_1}^{x_n} (g''(x) + h''(x))^2 dx \\ &= \int_{x_1}^{x_n} g''(x)^2 dx + \int_{x_1}^{x_n} h''(x)^2 dx + \int_{x_1}^{x_n} g''(x)h''(x) dx \end{aligned}$$

The third term is zero by integration by parts (Lebesgue-Stieltjes):

$$\int_{x_1}^{x_n} g''(x)h''(x)dx = \int_{x_1}^{x_n} g''(x)dh'(x) = \underbrace{g''h'|_{x_1}^{x_n}}_{=0} - \int h'dg''$$

Since  $g$  is a NCS, we have  $g''(x_n) = g''(x_1) = 0$  and  $g'''$  is a constant for any  $x$  within the interval  $(x_i, x_{i+1})$ ,  $i = 1, \dots, n-1$ . In other words,  $g'''$  is a piece-wise constant function on  $[a, b]$ .

$$\int h'dg'' = \sum_{i=1}^{n-1} g'''(x_i+) \int_{x_i}^{x_{i+1}} dh(x) = 0$$

where the last equality comes from the fact  $h(x_i) = 0$  for all  $i$ . Therefore,

$$\int_{x_1}^{x_n} f''(x)^2 dx \geq \int_{x_1}^{x_n} g''(x)^2 dx + \int_{x_1}^{x_n} h''(x)^2 dx \quad (10.10)$$

with "=" attained if and only if  $h''(x) = 0$  a.e. on  $[a, b]$ . However, since  $h \in W_2[a, b]$  and  $h(x_i) = 0$  for all  $i$ , we have (mean value theorem)

$$\frac{h'(z) - h'(y)}{z - y} = h''(x^*) = 0$$

for some  $x^*$  lies between  $(y, z)$  and  $y, z$  are two points within  $[a, b]$ . This means  $h'(x) = 0$  for a.e.  $x \in [a, b]$  and therefore we have  $h(x) = 0$  for a.e.  $x \in [a, b]$ . In other words any interpolant  $f$  will have a larger  $J(f)$  if  $f$  is not identical to  $g$  a.e.

□

## 10.4 The selection of knots position

### 10.4.1 Conventional methods

See these papers [21–23].

### 10.4.2 Bayesian methods

See these papers [32, 33].

### 10.4.3 Metaheuristic methods

See the paper written by the physicist Dr. Mohanty [24].



## References

- [1] George AF Seber and Alan J Lee. Linear regression analysis, volume 329. John Wiley & Sons, 2012.
- [2] Roger Penrose. A generalized inverse for matrices. In Mathematical proceedings of the Cambridge philosophical society, volume 51, pages 406–413. Cambridge University Press, 1955.
- [3] Calyampudi Radhakrishna Rao. Linear statistical inference and its applications, volume 2. Wiley New York, 1973.
- [4] Franklin A Graybill and George Marsaglia. Idempotent matrices and quadratic forms in the general linear hypothesis. The Annals of Mathematical Statistics, 28(3):678–686, 1957.
- [5] KS Banerjee. A note on idempotent matrices. The Annals of Mathematical Statistics, 35(2): 880–882, 1964.
- [6] R Dennis Cook and Sanford Weisberg. Diagnostics for heteroscedasticity in regression. Biometrika, 70(1):1–10, 1983.
- [7] Michael R Stoline and Hans K Ury. Tables of the studentized maximum modulus distribution and an application to multiple comparisons among means. Technometrics, 21(1):87–93, 1979.
- [8] Yosef Hochberg. Some generalizations of the t-method in simultaneous inference. Journal of multivariate analysis, 4(2):224–234, 1974.
- [9] Steven F Arnold. The theory of linear models and multivariate analysis. Wiley New York, 1981.
- [10] Mayuri Pandya, Krishnam Bhatt, and Paresh Andharia. Bayes estimation of two-phase linear regression model. Journal of Quality and Reliability Engineering, 2011.
- [11] David A Harville. Bayesian inference for variance components using only error contrasts. Biometrika, 61(2):383–385, 1974.
- [12] Lynn Roy LaMotte. A direct derivation of the reml likelihood function. Statistical Papers, 48 (2):321–327, 2007.
- [13] John D Storey, Wenzhong Xiao, Jeffrey T Leek, Ronald G Tompkins, and Ronald W Davis. Significance analysis of time course microarray experiments. Proceedings of the National Academy of Sciences, 102(36):12837–12842, 2005.

- [14] Cole Trapnell, Davide Cacchiarelli, Jonna Grimsby, Prapti Pokharel, Shuqiang Li, Michael Morse, Niall J Lennon, Kenneth J Livak, Tarjei S Mikkelsen, and John L Rinn. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. Nature biotechnology, 32(4):381–386, 2014.
- [15] Xu Ren and Pei-Fen Kuan. Negative binomial additive model for rna-seq data analysis. BMC bioinformatics, 21(1):1–15, 2020.
- [16] Koen Van den Berge, Hector Roux De Bezieux, Kelly Street, Wouter Saelens, Robrecht Cannoodt, Yvan Saeys, Sandrine Dudoit, and Lieven Clement. Trajectory-based differential expression analysis for single-cell sequencing data. Nature communications, 11(1):1–13, 2020.
- [17] Dongyuan Song and Jingyi Jessica Li. Pseudotimed: inference of differential gene expression along cell pseudotime with well-calibrated p-values from single-cell rna sequencing data. Genome biology, 22(1):1–25, 2021.
- [18] Rhonda Bacher, Ning Leng, Li-Fang Chu, Zijian Ni, James A Thomson, Christina Kendziorski, and Ron Stewart. Trendy: segmented regression analysis of expression dynamics in high-throughput ordered profiling experiments. BMC bioinformatics, 19(1):1–10, 2018.
- [19] Soumya Mohanty. Swarm intelligence methods for statistical regression. CRC Press, 2018.
- [20] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. An introduction to statistical learning, volume 112. Springer, 2013.
- [21] Hyungjun Park and Joo-Haeng Lee. B-spline curve fitting based on adaptive curve refinement using dominant points. Computer-Aided Design, 39(6):439–451, 2007.
- [22] Hongmei Kang, Falai Chen, Yusheng Li, Jiansong Deng, and Zhouwang Yang. Knot calculation for spline fitting via sparse optimization. Computer-Aided Design, 58:179–188, 2015.
- [23] Jiaqi Luo, Hongmei Kang, and Zhouwang Yang. Knot calculation for spline fitting based on the unimodality property. Computer Aided Geometric Design, 73:54–69, 2019.
- [24] Soumya D Mohanty and Ethan Fahnestock. Adaptive spline fitting with particle swarm optimization. Computational Statistics, 36(1):155–191, 2021.
- [25] Jerome H Friedman. The elements of statistical learning: Data mining, inference, and prediction. springer open, 2017.
- [26] Les Piegl and Wayne Tiller. The NURBS book. Springer Science & Business Media, 1996.
- [27] Carl De Boor. On calculating with b-splines. Journal of Approximation theory, 6(1):50–62, 1972.

- [28] Carl De Boor and Carl De Boor. A practical guide to splines, volume 27. springer-verlag New York, 1978.
- [29] Kenneth R Butterfield. The computation of all the derivatives of a b-spline basis. IMA Journal of Applied Mathematics, 17(1):15–25, 1976.
- [30] Isaac Jacob Schoenberg. On interpolation by spline functions and its minimal properties. In On Approximation Theory/Über Approximationstheorie, pages 109–129. Springer, 1964.
- [31] Peter J Green and Bernard W Silverman. Nonparametric regression and generalized linear models: a roughness penalty approach. Chapman and Hall/CRC, 2019.
- [32] Ilaria DiMatteo, Christopher R Genovese, and Robert E Kass. Bayesian curve-fitting with free-knot splines. Biometrika, 88(4):1055–1071, 2001.
- [33] Robert B Gramacy. Bayesian treed Gaussian process models. University of California, Santa Cruz, 2005.