

## Chapter 3. Linear regression for the full-rank model

The linear regression model is probably the most fundamental and widely used statistical model. Consider the following general linear model in matrix form

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}, \quad (1)$$

$\swarrow$   $p-1$  predictors.  
 $\nwarrow$  not r.v. but fixed design matrix.

where  $(Y, X^\top)$  is a pair of response and  $p$ -dimensional vector of covariates and  $\varepsilon$  is unobservable error term,  $\boldsymbol{\beta}_0$  is the true value of  $\boldsymbol{\beta}$ . The least squares (LS) and the least absolute deviation (LAD) are among the most widely-used criteria in statistical estimation for linear regression model. As a standard case, we consider  $\mathbf{X}$  is of full column rank, that is  $r(\mathbf{X}) = p$ .

for random design, we assume  $\mathbb{E}(\varepsilon|X) = 0$ .

### 3.1 Ordinary least squares estimation

there are many variant choice now.

For ordinary least squares estimation, it is commonly assumed that  $\mathbb{E}(\varepsilon) = \mathbf{0}$  and  $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$ , among which the mean-zero condition is an identifiability condition for the intercept component of  $\beta$ . The celebrated least squares estimate is to minimize

$Y = \alpha + X\beta + \varepsilon$  ,  $\mathbb{E}\varepsilon = 0$  .  $\text{Var}\varepsilon = \sigma^2 \mathbf{I}$   
 Note we do not assume  $\varepsilon \sim \mathcal{N}$

$$(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad \sum_{i=1}^n (y_i - x_i \beta)^2$$

$\mathbb{E}Y = X\beta$   
mean regression.

over  $\boldsymbol{\beta}$ . Simple calculations yields that

$$\begin{aligned} L(\boldsymbol{\beta}) &\equiv (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{Y}^\top \mathbf{Y} - 2\mathbf{Y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}. \end{aligned}$$

Then,

$$\frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

leads to the so-called *normal equation*

$$\begin{aligned} \mathbf{X}^\top \mathbf{Y} &= \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} \\ \Rightarrow \hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \end{aligned}$$

Compared with other existing methods, the LS is easy to implement and most popular, as its objective function  $L(\boldsymbol{\beta})$  is convex and the solution  $\hat{\boldsymbol{\beta}}$  is of a closed form.

Remark 1. Note that

$$\begin{aligned}
 L(\beta) &= (Y - X\beta)^\top (Y - X\beta) \\
 &= [Y - X\hat{\beta} + X(\hat{\beta} - \beta)]^\top [Y - X\hat{\beta} + X(\hat{\beta} - \beta)] \\
 &= (Y - X\hat{\beta})^\top (Y - X\hat{\beta}) + (\hat{\beta} - \beta)^\top X^\top X (\hat{\beta} - \beta) + 2(Y - X\hat{\beta})^\top X (\hat{\beta} - \beta).
 \end{aligned}$$

But

$$\begin{aligned}
 (Y - X\hat{\beta})^\top X &= (Y - X(X^\top X)^{-1}X^\top Y)^\top X \\
 &= 0.
 \end{aligned}$$

Then,

$$\begin{aligned}
 &(Y - X\beta)^\top (Y - X\beta) \\
 &= (Y - X\hat{\beta})^\top (Y - X\hat{\beta}) + (\hat{\beta} - \beta)^\top X^\top X (\hat{\beta} - \beta) \\
 &\geq 0,
 \end{aligned}$$

which achieves its minimum when  $\beta = \hat{\beta}$ .

Therefore,  $\hat{\beta} = (X^\top X)^{-1}X^\top Y$  is the least squares estimate of  $\beta_0$ .

### 3.1.1 Properties of the least squares estimate.

Given the least squares estimate, we define the vector of residuals as  $\hat{\varepsilon} = Y - X\hat{\beta}$ . Hence

$$\begin{aligned}
 \hat{\varepsilon} &= Y - X(X^\top X)^{-1}X^\top Y \\
 &= [I - X(X^\top X)^{-1}X^\top]Y \\
 &= [I - H]Y
 \end{aligned}$$

projection matrix that projects  $Y$   
onto the column space of  $X$ .

↳ projection?

where  $H = X(X^\top X)^{-1}X^\top$  is the so-called hat matrix of order  $n \times n$ . To obtain a fitted value of  $Y$ , we plug in  $\hat{\beta}$  and get

$$\hat{Y} = X\hat{\beta} = X(X^\top X)^{-1}X^\top Y = HY.$$

There are a number of properties here:

Properties:

1. The hat matrix  $H$  is symmetric idempotent;

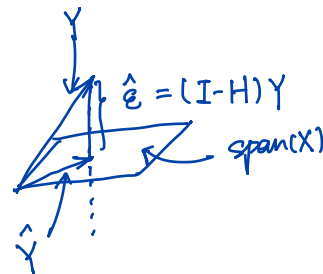
**Proof:** Note that  $(X^T X)$  is symmetric and  $(X^T X)^{-1}$  is also symmetric. Then,

$$H = X(X^T X)^{-1}X^T = H^T, \quad HH = X(X^T X)^{-1}X^T X(X^T X)^{-1}X^T = H.$$

2.  $X^T \hat{\varepsilon} = 0$ ; (This holds because of  $X^T H = X^T$ ,  $HX = X$  and  $X^T(I - H) = 0$ ,  $(I - H)X = 0$ .)

**Proof:** Since  
*projection matrix that project a vector onto the orthogonal complement.*  
 then,

$$X^T H = X^T X(X^T X)^{-1}X^T = X^T, \\
HX = X(X^T X)^{-1}X^T X = X,$$



$$X^T(I - H) = X^T - X^T H = X^T - X^T = 0, \\
(I - H)X = X - HX = X - X = 0.$$

Clearly,

$$X^T \hat{\varepsilon} = X^T(I - H)Y = 0.$$

3.  $\hat{Y}^T \hat{\varepsilon} = 0$ ;

**Proof:** Write

$$(HY)^T(I - H)Y = Y^T H^T(I - H)Y = Y^T H(I - H)Y \\
= Y^T 0Y = 0.$$

4.  $I - H$  is symmetric idempotent;

5.  $E(\hat{\beta}) = \beta_0$  (unbiased estimate);

**Proof:**

$$E(\hat{\beta}) = E((X^T X)^{-1}X^T Y) = (X^T X)^{-1}X^T E(Y) = (X^T X)^{-1}X^T X \beta_0 = \beta_0.$$

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

**Proof:**

$$Cov(\hat{\beta}) = Cov(\overset{A}{(X^T X)^{-1}X^T}Y) = \overset{A}{(X^T X)^{-1}X^T} \overset{\sigma^2 I}{Cov(Y)} \overset{A}{X(X^T X)^{-1}} \\
= \sigma^2 (X^T X)^{-1}X^T I X(X^T X)^{-1} = \sigma^2 (X^T X)^{-1}.$$

7.  $tr(\mathbf{I}_n - \mathbf{H}) = n - p$ ;

**Proof:** Note that

$$tr(\mathbf{H}) = \boxed{tr(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = tr(\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}) = tr(\mathbf{I}_p) = p.}$$

Then,

*Scalar.*

$$tr(\mathbf{I}_n - \mathbf{H}) = tr(\mathbf{I}_n) - tr(\mathbf{H}) = n - p.$$

8.  $\hat{\varepsilon}^\top \hat{\varepsilon} = tr(\mathbf{Y}\mathbf{Y}^\top(\mathbf{I} - \mathbf{H}))$ ;

**Proof:** We can easily show that

$$\begin{aligned} \hat{\varepsilon}^\top \hat{\varepsilon} &= \mathbf{Y}^\top (\mathbf{I} - \mathbf{H})^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y} = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y} \\ &= tr(\mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}) = tr(\mathbf{Y}\mathbf{Y}^\top (\mathbf{I} - \mathbf{H})). \end{aligned}$$

9.  $E(\mathbf{Y}\mathbf{Y}^\top) = \sigma^2 \mathbf{I} + \mathbf{X}\beta\beta^\top \mathbf{X}^\top$ ;

10.  $\hat{\varepsilon}^\top \hat{\varepsilon} / (n - p)$  is an unbiased estimate of  $\sigma^2$ , that is  
*RSS.*

$$E\left(\frac{\hat{\varepsilon}^\top \hat{\varepsilon}}{n - p}\right) = \sigma^2.$$

**Proof:** Write

$$\begin{aligned} E(\hat{\varepsilon}^\top \hat{\varepsilon}) &= E(\mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H}) \mathbf{Y}) = E(\mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}) \\ &= tr((\mathbf{I} - \mathbf{H}) \Sigma) + \beta^\top \mathbf{X}^\top (\mathbf{I} - \mathbf{H}) \mathbf{X} \beta \\ &= \sigma^2 tr(\mathbf{I} - \mathbf{H}) = \sigma^2 (n - p). \end{aligned}$$

*0. I-H ⊥ X.*

Thus,  $E\left(\frac{\hat{\varepsilon}^\top \hat{\varepsilon}}{n - p}\right) = \sigma^2$ .  *$\min_{f \in \mathcal{B}} E(Y - f(X))^2 \Rightarrow \hat{f} = E(Y|X)$ .*

*Remark 2.* Note that in this course, we mostly consider fixed design, that is the covariate  $X$  is fixed and deterministic. For random design, the least square estimation is still valid and its theoretical properties can be established without further difficulties.

### 3.2 The weighted least square estimation.

For a general case that  $Cov(\varepsilon) = \Sigma$  and  $\Sigma$  is known, the weighted least squares will be used

*↑  
Bootstrap to estimate the  $\Sigma$  of  $\varepsilon$ .*

to estimate  $\beta$  in model (1). Note that  $\Sigma \neq \mathbf{I}$  in general but is positive definite. Recall that the ordinary least squares is to minimize  $(\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)$  and  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ . The weighted least squares (WLS) or generalized least squares (GLS) estimator is defined as the minimizer of

$$(\mathbf{Y} - \mathbf{X}\beta)^\top \Sigma^{-1} (\mathbf{Y} - \mathbf{X}\beta)$$

over  $\beta$ .

Similar to section 2.1, we let

$$\begin{aligned} S(\beta) &= (\mathbf{Y} - \mathbf{X}\beta)^\top \Sigma^{-1} (\mathbf{Y} - \mathbf{X}\beta) \\ &= \mathbf{Y}^\top \Sigma^{-1} \mathbf{Y} - 2\mathbf{Y}^\top \Sigma^{-1} \mathbf{X}\beta + \beta^\top \mathbf{X}^\top \Sigma^{-1} \mathbf{X}\beta. \end{aligned}$$

Then,

$$\begin{aligned} \frac{\partial S(\beta)}{\partial \beta} &= -2\mathbf{X}^\top \Sigma^{-1} \mathbf{Y} + 2\mathbf{X}^\top \Sigma^{-1} \mathbf{X}\beta \\ &= \mathbf{0} \\ \Rightarrow \mathbf{X}^\top \Sigma^{-1} \mathbf{Y} &= \mathbf{X}^\top \Sigma^{-1} \mathbf{X}\beta \\ \Rightarrow \tilde{\beta} &= (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Sigma^{-1} \mathbf{Y}. \end{aligned}$$

Note that  $E(\tilde{\beta}) = \beta_0$  and  $Cov(\tilde{\beta}) = (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1}$ .

*Remark 3.* When  $\Sigma = \sigma^2 \mathbf{I}$ , the WLS or GLS reduces to the OLS.

*Remark 4.* We provide another aspect to motivate the WLS. Since  $\Sigma$  is positive definite,  $\Sigma^{-1/2}$  exists such that  $\Sigma^{-1/2} \Sigma^{-1/2} = \Sigma^{-1}$ . Thus,

$$\boxed{\Sigma^{-\frac{1}{2}} \mathbf{Y}} = \Sigma^{-\frac{1}{2}} \mathbf{X}\beta + \Sigma^{-\frac{1}{2}} \epsilon.$$

Now  $E(\Sigma^{-\frac{1}{2}} \epsilon) = \mathbf{0}$  and  $\boxed{Cov(\Sigma^{-\frac{1}{2}} \epsilon) = \mathbf{I}_n}$  satisfy the conditions of the ordinary least squares.

Thereby,

$$\begin{aligned} \tilde{\beta} &= \left\{ (\Sigma^{-\frac{1}{2}} \mathbf{X})^\top (\Sigma^{-\frac{1}{2}} \mathbf{X}) \right\}^{-1} (\Sigma^{-\frac{1}{2}} \mathbf{X})^\top \Sigma^{-\frac{1}{2}} \mathbf{Y} \\ &= (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Sigma^{-1} \mathbf{Y}. \end{aligned}$$

*actually it is rotation and scaling.*

NEED TO KNOW HOW TO PROOF B.L.U.E.

### 3.3 The Best linear unbiased estimator (b.l.u.e. or BLUE) (Gauss-Markov Theorem)

Let  $\mathbf{t} \in \mathbb{R}^p$  be a vector. We consider the problem of finding the b.l.u.e. of  $\mathbf{t}^\top \boldsymbol{\beta}$ . Let  $\boldsymbol{\lambda}^\top \mathbf{Y}$  be a linear function of the observations and an estimator of  $\mathbf{t}^\top \boldsymbol{\beta}$ . To find the BLUE of  $\mathbf{t}^\top \boldsymbol{\beta}$  is to determine  $\boldsymbol{\lambda}$  such that  $\boldsymbol{\lambda}^\top \mathbf{Y}$  is unbiased for  $\mathbf{t}^\top \boldsymbol{\beta}$  and has minimum variance among all the linear unbiased estimates. To this end,

1. First, if  $\boldsymbol{\lambda}^\top \mathbf{Y}$  is an unbiased estimator of  $\mathbf{t}^\top \boldsymbol{\beta}$ ,  $E(\boldsymbol{\lambda}^\top \mathbf{Y}) = \mathbf{t}^\top \boldsymbol{\beta}$ . But  $E(\boldsymbol{\lambda}^\top \mathbf{Y}) = \boldsymbol{\lambda}^\top E(\mathbf{Y}) = \boldsymbol{\lambda}^\top \mathbf{X}\boldsymbol{\beta}$  according to model (1). Hence,

$$\boldsymbol{\lambda}^\top \mathbf{X}\boldsymbol{\beta} = \mathbf{t}^\top \boldsymbol{\beta}$$

which is true for all  $\boldsymbol{\beta}$ . Thus,  $\boldsymbol{\lambda}^\top \mathbf{X} = \mathbf{t}^\top$ . *Unbiased. serves as constrain*

2. Second, we need to find the linear unbiased estimator of  $\mathbf{t}^\top \boldsymbol{\beta}$  which has minimum variance. Note that

$$Var(\boldsymbol{\lambda}^\top \mathbf{Y}) = \boldsymbol{\lambda}^\top \boldsymbol{\Sigma} \boldsymbol{\lambda}.$$

Using  $2\boldsymbol{\theta}$  as a vector of Lagrange multipliers, we need to minimize

$$W(\boldsymbol{\lambda}, \boldsymbol{\theta}) = \boldsymbol{\lambda}^\top \boldsymbol{\Sigma} \boldsymbol{\lambda} - 2\boldsymbol{\theta}^\top (\mathbf{X}^\top \boldsymbol{\lambda} - \mathbf{t}),$$

where  $\mathbf{X}^\top \boldsymbol{\lambda} = \mathbf{t}$  is the unbiasedness condition. Thus,

$$\left. \begin{aligned} \frac{\partial W(\boldsymbol{\lambda}, \boldsymbol{\theta})}{\partial \boldsymbol{\lambda}} &= 2\boldsymbol{\Sigma} \boldsymbol{\lambda} - 2\mathbf{X}\boldsymbol{\theta} = 0, \\ \frac{\partial W(\boldsymbol{\lambda}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= 2\mathbf{X}^\top \boldsymbol{\lambda} - 2\mathbf{t} = 0. \end{aligned} \right\} ?$$

Solving the above two equations for  $\boldsymbol{\lambda}$  and  $\boldsymbol{\theta}$ , we have

$$\boldsymbol{\lambda}^\top = \mathbf{t}^\top (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1}.$$

*$\boldsymbol{\lambda}^\top = \mathbf{t}^\top \hat{\boldsymbol{\beta}}_{WLS}$*

Therefore, the BLUE of  $\mathbf{t}^\top \boldsymbol{\beta}$  is

$$\boldsymbol{\lambda}^\top \mathbf{Y} = \mathbf{t}^\top (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{Y},$$

with variance

$$\begin{aligned} Var(\boldsymbol{\lambda}^\top \mathbf{Y}) &= \mathbf{t}^\top (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\Sigma}) (\boldsymbol{\Sigma}^{-1}) \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{t} \\ &= \mathbf{t}^\top (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{t}. \end{aligned}$$

*Remark 5.* In a special case that  $\Sigma = \sigma^2 \mathbf{I}$ , the BLUE of  $\mathbf{t}^\top \beta$  is

$$\mathbf{t}^\top (\mathbf{X}^\top (\mathbf{I} \sigma^2)^{-1} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I} \sigma^2)^{-1} \mathbf{Y} = \mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y},$$

with variance

$$\mathbf{t}^\top (\mathbf{X}^\top (\mathbf{I} \sigma^2)^{-1} \mathbf{X})^{-1} \mathbf{t} = \sigma^2 \mathbf{t}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{t}.$$

*Remark 6.* By letting  $\mathbf{t}^\top$  be, in turn, each row of  $\mathbf{I}_k$ , we can easily obtain the BLUE of  $\beta = \tilde{\beta} = (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Sigma^{-1} \mathbf{Y}$ , which is precisely the weighted least square estimate or generalized least square estimate.

*Remark 7.* When  $\Sigma = \sigma^2 \mathbf{I}$ , the BLUE of  $\beta$  is  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ .

**In summary, the least square estimate of  $\beta_0$  in (1) is the best linear unbiased estimate.**

**THEOREM 1.**  $W = \lambda^\top \Sigma \lambda$  is minimum if

$$\lambda^\top = \mathbf{t}^\top (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Sigma^{-1}$$

subject to the constraint that

$$\mathbf{X}^\top \lambda = \mathbf{t}.$$

*Proof.* Let  $\lambda_1^\top = \mathbf{t}^\top (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Sigma^{-1}$ . Let  $\lambda_2$  be another vector that is different from  $\lambda$  but also satisfies  $\mathbf{X}^\top \lambda_2 = \mathbf{t}$ . Then,

$$\begin{aligned} W^\top &= \lambda_2^\top \Sigma \lambda_2 \\ &= [(\lambda_2 - \lambda_1) + \lambda_1]^\top \Sigma [(\lambda_2 - \lambda_1) + \lambda_1] \\ &= (\lambda_2 - \lambda_1)^\top \Sigma (\lambda_2 - \lambda_1) + \lambda_1^\top \Sigma \lambda_1 + (\lambda_2 - \lambda_1)^\top \Sigma \lambda_1 + \lambda_1^\top \Sigma (\lambda_2 - \lambda_1). \end{aligned}$$

Nevertheless,

$$\begin{aligned} (\lambda_2 - \lambda_1)^\top \Sigma \lambda_1 &= (\lambda_2 - \lambda_1)^\top \Sigma [\Sigma^{-1} \mathbf{X} (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{t}] \\ &= (\lambda_2 - \lambda_1)^\top \mathbf{X} (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{t} \\ &= 0 \text{ (this is because } \lambda_1^\top \mathbf{X} = \mathbf{t}^\top \text{ and } \lambda_2^\top \mathbf{X} = \mathbf{t}^\top \text{)}. \end{aligned}$$

Also,

$$\lambda_1^\top \Sigma (\lambda_2 - \lambda_1) = (\lambda_2 - \lambda_1)^\top \Sigma \lambda_1 = 0.$$

As a result,

$$W^\top = (\lambda_2 - \lambda_1)^\top \Sigma (\lambda_2 - \lambda_1) + \lambda_1^\top \Sigma \lambda_1.$$

which is minimized if  $\lambda_2 = \lambda_1$ . The proof is complete.  $\square$

$$y_i = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_{p-1} x_{p-1,i} + \varepsilon_i$$

$\beta_0, \dots, \beta_{p-1}$  fixed effect  
 $\varepsilon_i$  random

### 3.4 Least squares theory when the parameters are random variables (random-effect model)

In this section, we assume that the parameters of the regression models are random variables with a known mean and variance. Consider the linear model

$$Y = Xb + e, \quad (2)$$

$(= \sum \beta_j X_{j,i} + e_i)$

where  $(Y_i, b_i, e_i), i = 1, \dots, n$  are independent and identically distributed (i.i.d) copies of  $(Y, b, e)$ , and  $E(b) = \theta$  and  $Cov(b) = F$ ,  $\theta$  is a  $k$ -dimensional vector and  $F$  is a  $k \times k$  positive definite matrix. Also assume that

*Notation : estimate a fixed value.  
predict a random variable*

*Assumptions in random effect linear model.*

$$E(e|b) = 0, \quad Cov(e|b) = V.$$

We then show how to find the best linear estimator (predictor) of a random variable  $p^\top b$ , where  $p \in \mathbb{R}^k$  is a given vector. The following formulae connect the conditional and unconditional means and variances.

$$\begin{aligned} E(Y) &= E(E(Y|e)), \\ Var(Y) &= E\{Var(Y|b)\} + Var\{E(Y|b)\} = V + XF^\top X^\top, \\ Cov(Y, p^\top b) &= E\{Cov(Y, p^\top b|b)\} + Cov[E(Y|b), p^\top b] = XFp. \end{aligned} \quad (3)$$

Students need to show the above formula by themselves as basic exercises on conditional expectation. The third equation above is by the **law of total covariance**, that is,

$$Cov(X, Y) = E[Cov(X, Y|Z)] + Cov(E(X|Z), E(Y|Z)).$$

The objective is to determine a linear function  $a + L^\top Y$  such that

*not sure  $L^\top Y$  is unbiased*

$$\begin{aligned} E(p^\top b - a - L^\top Y) &= E(p^\top b) - a - E(L^\top Y) = p^\top \theta - a - E(L^\top Y) \\ E(L^\top Y) &= E(L^\top (Xb + e)) \\ &= E[E(L^\top (Xb + e)|b)] \\ &= E[L^\top Xb + E(L^\top e|b)] \\ &= L^\top X\theta \\ 0 &= E(p^\top b - a - L^\top Y) = p^\top \theta - a - L^\top X\theta \Rightarrow a = (p^\top - L^\top X)\theta. \end{aligned} \quad (4)$$



$$\text{Var}(A - a - B) = \text{Var}A + \text{Var}B - 2\text{Cov}(A, B)$$

$$\begin{aligned}\text{Var}Y &= \mathbb{E}(\text{Var}(Y|b)) + \text{Var}(\mathbb{E}(Y|b)) \\ &= \mathbb{E}(\text{Var}(Xb + \varepsilon|b)) + \text{Var}(\mathbb{E}(Xb + \varepsilon|b)) \\ &= \mathbb{E}(\text{Var}(\varepsilon|b)) + \text{Var}(Xb) \\ &= V + XFXT^\top\end{aligned}$$

$$\Rightarrow \text{Var}(L^\top Y) = L^\top (V + XFXT^\top) L \equiv \text{Var}(\mathbf{p}^\top \mathbf{b} - a - \mathbf{L}^\top \mathbf{Y}) \text{ achieves its minimum.} \quad (5)$$

$$\begin{aligned}\text{then } v &= \text{Var}(\mathbf{p}^\top \mathbf{b}) + \text{Var}(\mathbf{L}^\top \mathbf{Y}) - 2\text{Cov}(\mathbf{p}^\top \mathbf{b}, \mathbf{L}^\top \mathbf{Y}) \\ &= \mathbf{p}^\top \mathbf{F} \mathbf{p} + \mathbf{L}^\top (\mathbf{X} \mathbf{F} \mathbf{X}^\top + \mathbf{V}) \mathbf{L} - 2 \mathbf{L}^\top \mathbf{X} \mathbf{F} \mathbf{p}.\end{aligned}$$

**THEOREM 2.** The optimum estimator/predictor that satisfies (4) and (5) takes the form

$$\mathbf{p}^\top \hat{\mathbf{b}} = \mathbf{p}^\top \boldsymbol{\theta} + \mathbf{p}^\top \mathbf{F} \mathbf{X}^\top (\mathbf{V} + \mathbf{X} \mathbf{F} \mathbf{X}^\top)^{-1} (\mathbf{Y} - \mathbf{X} \boldsymbol{\theta}) \quad (6)$$

$$= \mathbf{p}^\top \boldsymbol{\theta} + \mathbf{p}^\top (\mathbf{F}^{-1} + \mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X} \boldsymbol{\theta}). \quad (7)$$

**Proof:** The expectation in (4) yields

$$a = (\mathbf{p}^\top - \mathbf{L}^\top \mathbf{X}) \boldsymbol{\theta}. \quad (8)$$

Employing the three formula in (3), the quantity to be minimized in (5) is

$$v = \underline{\mathbf{p}^\top \mathbf{F} \mathbf{p} + \mathbf{L}^\top (\mathbf{X} \mathbf{F} \mathbf{X}^\top + \mathbf{V}) \mathbf{L} - 2 \mathbf{L}^\top \mathbf{X} \mathbf{F} \mathbf{p}}.$$

Then, differentiating  $v$  with respect to  $\mathbf{L}$  and setting the results equal to zero, we obtain

$$(\mathbf{X} \mathbf{F} \mathbf{X}^\top + \mathbf{V}) \mathbf{L} = \mathbf{X} \mathbf{F} \mathbf{p}$$

and the optimizing  $\mathbf{L}$  is

$$\underline{\mathbf{L}^* = (\mathbf{X} \mathbf{F} \mathbf{X}^\top + \mathbf{V})^{-1} \mathbf{X} \mathbf{F} \mathbf{p}}. \quad (9)$$

Substitution of (8) and (9) into  $a + \mathbf{L}^\top \mathbf{Y}$  yields (6). The equivalence of the two expressions in (7) is established by using the following Woodbury (1950) matrix identity

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{C}^{-1} + \mathbf{D} \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{D} \mathbf{A}^{-1},$$

where  $\mathbf{A} = \mathbf{V}$ ,  $\mathbf{B} = \mathbf{X}$ ,  $\mathbf{C} = \mathbf{F}$  and  $\mathbf{D} = \mathbf{X}^\top$ . The proof is complete.

Substitution into (9) gives the minimum variance

$$\begin{aligned}v_{\min} &= \mathbf{p}^\top \mathbf{F} \mathbf{p} - \mathbf{p}^\top \mathbf{F} \mathbf{X}^\top (\mathbf{X} \mathbf{F} \mathbf{X}^\top + \mathbf{V})^{-1} \mathbf{X} \mathbf{F} \mathbf{p} \\ &= \mathbf{p}^\top (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{p} - (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} (\mathbf{F} + (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1})^{-1} (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{p}.\end{aligned}$$

Notice that  $v_{min}$  is less than the variance of the least-square estimator.

*Remark 8.* When  $\mathbf{F} = \sigma^2 \mathbf{G}^{-1}$ ,  $\mathbf{V} = \sigma \mathbf{I}$  and  $\boldsymbol{\theta} = \mathbf{0}$ , the estimator in (6) reduces to

$$\mathbf{p}^\top \hat{\mathbf{b}} = \mathbf{p}^\top (\mathbf{X}^\top \mathbf{X} + \mathbf{G})^{-1} \mathbf{X}^\top \mathbf{Y},$$

the *generalized ridge regression* estimator of C.R. Rao (1975). When  $\mathbf{G} = k\mathbf{I}$ , it reduces to the ridge regression estimator of Hoerl and Kennard (1970). We will introduce the ridge regression in details in later sections.