

# **Chapter 3**

## **Multiple Linear Regression**

# Chapter 3b

Properties of least-squares estimators (pages 3-6)

Gauss Markov theorem (page 7)

Residual sum of squares and estimation of  $\sigma^2$  (pages 8-9)

Example – Cholesterol age data (pages 10-11)

An example to show that  $MS_{Res}$  is model dependent

Example – Delivery time data (page 12-15)

Inadequacy of scatter diagrams in multiple regression (page 16)

Maximum-likelihood estimation (pages 17-18)

Distribution theory – An important theorem (pages 19-20)

## Properties of least-squares estimators

$\underset{\text{identity}}{\overset{\text{nxn}}{\uparrow}}$   
matrix

1. Multiple linear regression model

$$\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times (k+1)} \boldsymbol{\beta}_{(k+1) \times 1} + \boldsymbol{\epsilon}_{n \times 1} \text{ where } E(\boldsymbol{\epsilon}) = \mathbf{0} \text{ and } Var(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$$

Least-squares normal equations

$$\mathbf{X}' \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}' \mathbf{y}$$

Least-squares estimators

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}$$

2. Note: Let  $\mathbf{A}$  be a  $k \times k$  matrix of constants and  $\mathbf{y}$  be a  $k \times 1$  random vector with mean  $E(\mathbf{y}) = \boldsymbol{\mu}$  and non-singular variance-covariance matrix  $Var(\mathbf{y}) = \mathbf{V}$ .

Then

(a)  $E(\mathbf{A}\mathbf{y}) = \mathbf{A}\boldsymbol{\mu}$

$$Var(a\mathbf{y}) = a'Var(\mathbf{y})a$$

(b)  $Var(\mathbf{A}\mathbf{y}) = \mathbf{A}\mathbf{V}\mathbf{A}'$

3.  $\hat{\beta}$  is an unbiased estimator of  $\beta$

$$\begin{aligned} E(\hat{\beta}) &= E[(X'X)^{-1}X'y] \\ &= E[(X'X)^{-1}X'(X\beta + \epsilon)] \\ &= E[(X'X)^{-1}X'X\beta + (X'X)^{-1}X'\epsilon] \\ &= \beta + (X'X)^{-1}\underbrace{X'E[\epsilon]}_{\substack{\text{assumption} \\ \text{that } E[\epsilon] = 0}} \end{aligned}$$

#### 4. Variance-covariance matrix of $\hat{\beta}$

$$\begin{aligned}
 \text{Var}(\hat{\beta}) &= \text{Var}[(X'X)^{-1}X'y] \\
 &= (X'X)^{-1}X'\text{Var}(y)[(X'X)^{-1}X']' \\
 &= (X'X)^{-1}X'\underbrace{\sigma^2 I_n}_{\text{circled}} X(X'X)^{-1} \\
 &= \sigma^2 (X'X)^{-1}X'X(X'X)^{-1}
 \end{aligned}$$

$$\begin{aligned}
 &= \sigma^2 (X'X)^{-1} \\
 &= (X')' \left[ (X'X)^{-1} \right]' \\
 &= X \cdot (X'X)^{-1}
 \end{aligned}$$

$X'X$  is symmetric  
 $\Rightarrow (X'X)^{-1}$   
 is also symmetric

$$Var \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \vdots \\ \hat{\beta}_k \end{bmatrix} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 \mathbf{C} = \sigma^2 \begin{bmatrix} C_{00} & C_{01} & C_{02} & \dots & C_{0,k} \\ C_{10} & C_{11} & C_{12} & \dots & C_{1,k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \dots & \vdots \\ C_{k0} & C_{k1} & C_{k2} & \dots & C_{k,k} \end{bmatrix}$$

$\text{Var}(\hat{\beta}_0) \leftarrow$   
  
 $\text{Var}(\hat{\beta}_1) \rightarrow$

$$Var(\hat{\beta}_j) = \sigma^2 C_{jj}$$

$$Cov(\hat{\beta}_i, \hat{\beta}_j) = \sigma^2 C_{ij}$$

## Gauss-Markov theorem

$\hat{\beta}$  is the best linear unbiased estimator of  $\beta$

$\hat{\beta}$  is also the maximum-likelihood estimator of  $\beta$ .

## Residual sum of squares and estimation of $\sigma^2$

Let  $SS_{Res}$  be the residual sum of squares.

$$\begin{aligned} SS_{Res} &\equiv \sum_{i=1}^n (y_i - \hat{y}_i)^2 & E(\varepsilon) = 0 \\ &= \sum_{i=1}^n e_i^2 & \text{Var}(\varepsilon) = \sigma^2 I_n \\ &= \mathbf{e}'\mathbf{e} \\ &= (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\ &= (\mathbf{y}' - \hat{\boldsymbol{\beta}}'\mathbf{X}')(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\ &= \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} - \mathbf{y}'\mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{y}'\mathbf{y} - 2\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} & \because \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y} \\ SS_{Res} &= \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} \end{aligned}$$

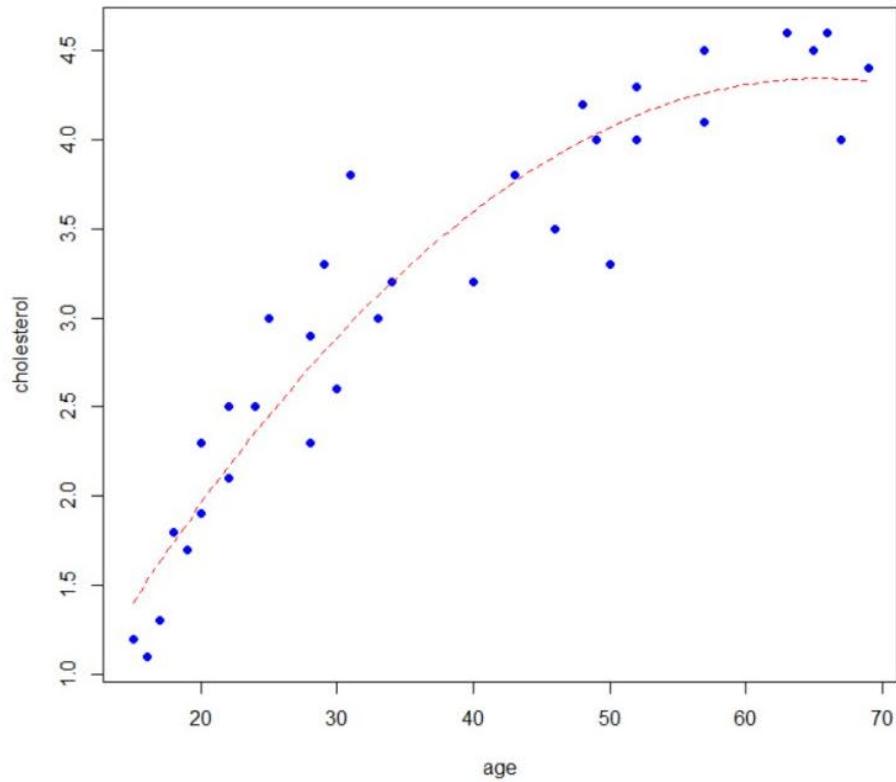
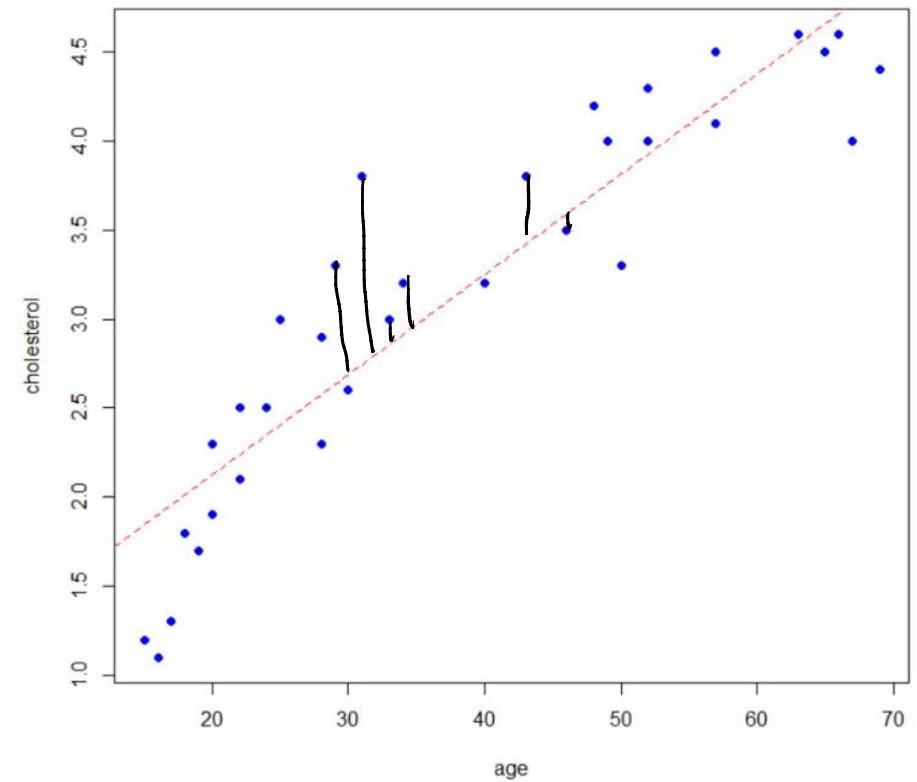
$$SS_{Res} = \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y}$$

$$\hat{\sigma}^2 = MS_{Res} \equiv \frac{SS_{Res}}{n - p}$$

Linear model theory:  $E(MS_{Res}) = \sigma^2$

$MS_{Res}$  is model-dependent because  $MS_{Res}$  changes with the model.

## Cholesterol age data



```
48 # Compare straight line and quadratic models  
49 # in terms of MS_Res  
50 anova(lm(cholesterol~age))  
51 anova(lm(cholesterol~age + I(age^2)))
```

48:1 (Top Level) ↴

Console Background Jobs ✎

R 3.4.1 · ↶  
> # Compare straight line and quadratic models

> # in terms of MS\_Res

> anova(lm(cholesterol~age))

Analysis of Variance Table

Response: cholesterol

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
age	1	30.7370	30.7370	161.52	7.856e-14 ***
Residuals	31	5.8993	0.1903		

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

> anova(lm(cholesterol~age + I(age^2)))

Analysis of Variance Table

Response: cholesterol

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
age	1	30.7370	30.7370	262.39	2.234e-16 ***
I(age^2)	1	2.3851	2.3851	20.36	9.198e-05 ***
Residuals	30	3.5143	0.1171		

$$MS_{REJ} = \sigma^2$$

**Residual mean square is smaller for a better fitted model**

## Delivery Time Data

Column of response

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{y}' = \text{transpose of } \mathbf{y} \\ = (y_1, y_2, \dots, y_n)$$

$$\mathbf{y}'\mathbf{y} = \sum_{i=1}^{25} y_i^2 = 18,310.6290$$

$\hat{\beta}' \mathbf{X}' \mathbf{y}$  = [2.34123115 1.61590721 0.01438483]  $\begin{bmatrix} 559.60 \\ 7,375.44 \\ 337,072.00 \end{bmatrix}$

= 18,076.90304

$$SS_{Res} = \mathbf{y}'\mathbf{y} - \hat{\beta}' \mathbf{X}' \mathbf{y}$$
$$= 18,310.6290 - 18,076.9030 = 233.7260$$

$$MS_{Res} \hat{\sigma}^2 = \frac{SS_{Res}}{n-p} = \frac{233.7260}{25-3} = 10.6239$$

$\nwarrow \uparrow$   
K+1

```
34 # matrices and finding betas
35 one <- rep(1,length(time))
36 X <- array(c(one,case,dist), dim=c(length(time),3))
37 y <- array(time, dim=c(length(time),1))
38 XPX <- t(X) %*% X
39 XPX
40 XPy <- t(X) %*% y
41 XPy
42 betahat <- solve(XPX, XPy)
43 betahat

$$x' \hat{x} = x'y \quad \text{use solve function}$$

44
45 # estimator of sigma^2
46 bxy <- t(betahat) %*% t(X) %*% y
47 bxy
48 SSRes <- t(y) %*% y - bxy
49 SSRes
50 sigma2hat <- SSRes/(length(time)-3)
51 sigma2hat
52
53 # alternative method
54 anova(fitted.model)
```

```
> # matrices and finding betas
> one <- rep(1,length(time))
> X <- array(c(one,case,dist), dim=c(length(time),3))
> y <- array(time, dim=c(length(time),1))
> XPX <- t(X) %*% X
> XPX
      [,1]   [,2]   [,3]
[1,]    25   219 10232
[2,]   219  3055 133899
[3,] 10232 133899 6725688
> XPy <- t(X) %*% y
> XPy
      [,1]
[1,] 559.60
[2,] 7375.44
[3,] 337071.69
> betahat <- solve(XPX, XPy)
> betahat
      [,1]
[1,] 2.34123115
[2,] 1.61590721
[3,] 0.01438483
```

```

> # estimator of sigma^2
> bxy <- t(betahat) %*% t(X) %*% y
> bxy
      [,1]
[1,] 18076.9
> SSRes <- t(y) %*% y - bxy
> SSRes
      [,1]
[1,] 233.7317
> sigma2hat <- SSRes/(length(time)-3)
> sigma2hat
      [,1]  $\hat{\sigma}^2$ 
[1,] 10.62417
> # alternative method
> anova(fitted.model)
Analysis of Variance Table

```

Response: time

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
case	1	5382.4	5382.4	506.619	< 2.2e-16	***
dist	1	168.4	168.4	15.851	0.0006312	***
Residuals	22	233.7	10.6 $\hat{\sigma}^2$			

SSRes

## Inadequacy of Scatter Diagrams in Multiple Regression

These data were generated from the equation

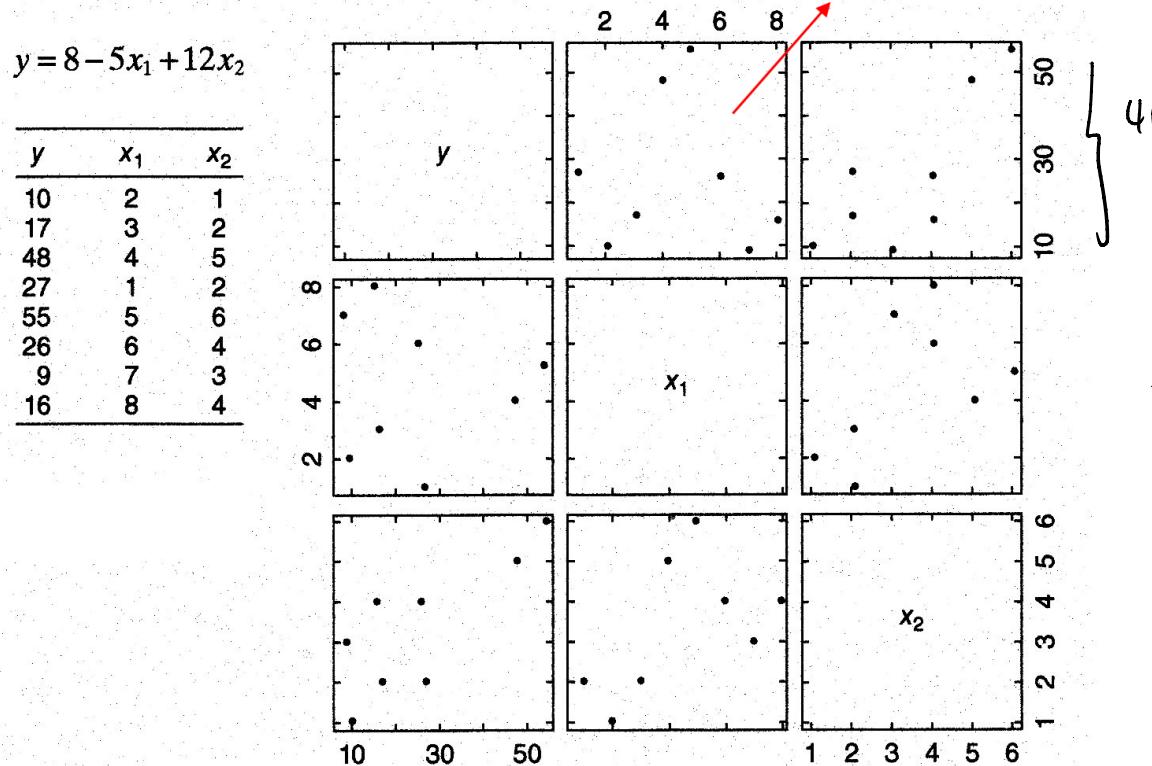
$$y = 8 - 5x_1 + 12x_2$$

The matrix of scatterplots is shown. The  $y$ -versus- $x_1$ , plot does not exhibit any apparent relationship between the two variables. The  $y$ -versus- $x_2$  plot indicates that a linear relationship exists, with a slope of approximately 8. Note that both scatter diagrams convey erroneous information. Since in this data set there are two pairs of points that have the same  $x_2$  values ( $x_2 = 2$  and  $x_2 = 4$ ), we could measure the  $x_1$  effect at fixed  $x_2$  from both pairs. This gives,  $\hat{\beta}_1 = (17 - 27)/(3 - 1) = -5$  for  $x_2 = 2$  and  $\hat{\beta}_1 = (26 - 16)/(6 - 8) = -5$  for  $x_2 = 4$  the correct results.

simulated  
data

$$y = 8 - 5x_1 + 12x_2$$

$y$	$x_1$	$x_2$
10	2	1
17	3	2
48	4	5
27	1	2
55	5	6
26	6	4
9	7	3
16	8	4



## Maximum-Likelihood Estimation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \Rightarrow \mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

$\boldsymbol{\varepsilon}$  is distributed as  $N(\mathbf{0}, \sigma^2 \mathbf{I})$

$$f(\varepsilon_i) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2} \varepsilon_i^2\right)$$

$$\begin{aligned} \text{Var}(\boldsymbol{\varepsilon}) &= \sigma^2 \underbrace{\begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}}_{\text{identity matrix}} \\ &= \begin{pmatrix} \sigma^2 & & & \\ & \sigma^2 & & \\ & & \ddots & \\ & & & \sigma^2 \end{pmatrix} \Leftrightarrow \text{Var}(\varepsilon_i) = \sigma^2 \end{aligned}$$

Var( $\varepsilon_1$ ) Var( $\varepsilon_2$ )

Cov( $\varepsilon_i, \varepsilon_j$ ,  
 $i \neq j$ )

The likelihood function is the joint density of  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  or  $\prod_{i=1}^n f(\varepsilon_i)$

$$\begin{matrix} & \text{scalar} \\ L(\boxed{\boldsymbol{\varepsilon}, \boldsymbol{\beta}}, \sigma^2) &= \prod_{i=1}^n f(\varepsilon_i) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left(-\frac{1}{2\sigma^2} \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}\right) \\ & \text{matrix} \end{matrix}$$

$$L(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$$

$$\begin{array}{c} \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ \downarrow \\ \boldsymbol{\varepsilon} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \end{array}$$

$$\ln L(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \ln(2\pi) - n \ln(\sigma) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

**Maximize**  $\ln L(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \sigma^2)$  is the same as

**minimize**  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  for a fixed value of  $\sigma$ .

Least-squares function

Maximum-likelihood estimator of  $\boldsymbol{\beta}$  under normal errors  
is the same as the least-squares estimator  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$ .

The maximum-likelihood estimator of  $\sigma^2$  is

$$\tilde{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n}$$

$$\hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n-p}$$

$\nwarrow E(\hat{\sigma}^2) \neq \sigma^2$

$$E(\hat{\sigma}^2) = \hat{\sigma}_2^2$$

(unbiased  
↳ preferred)

## Distribution theory (graduate level course in multiple linear regression)

**Theorem DT:** Let  $\mathbf{A}$  be a  $k \times k$  matrix of constants, and  $\mathbf{y}$  be a  $k \times 1$  multivariate normal vector with mean  $\boldsymbol{\mu}$  and non-singular variance-covariance matrix  $\mathbf{V}$ ; thus

$$\mathbb{E}(\mathbf{y}) = \boldsymbol{\mu}, \quad \text{Var}(\mathbf{y}) = \mathbf{V} \\ \mathbf{y} \sim N(\boldsymbol{\mu}, \mathbf{V}).$$

Let  $U$  be the quadratic form defined by  $U = \mathbf{y}' \mathbf{A} \mathbf{y}$ .  $\forall \mathbf{y}, \mathbf{A}\mathbf{A} = \mathbf{A}$

1. If  $\mathbf{AV}$  or  $\mathbf{VA}$  is an idempotent matrix of rank  $p$ , then

$$U \sim \chi_{p,\lambda}^2$$

where  $\chi_{p,\lambda}^2$  is the chi-square distribution with  $p$  degrees of freedom and non-centrality parameter  $\lambda = \boldsymbol{\mu}' \mathbf{A} \boldsymbol{\mu}$ .

2. Let  $\mathbf{V} = \sigma^2 \mathbf{I}$ , which is a typical assumption. If  $\mathbf{A}$  is an idempotent matrix of rank  $p$ , then

$$U/\sigma^2 \sim \chi_{p,\lambda}^2 \Rightarrow \underbrace{\frac{(n-p) \text{MS}_{\text{Res}}}{\sigma^2}}_{\text{MS}_{\text{Res}}} \sim \chi_{n-p}^2$$

where  $\lambda = \boldsymbol{\mu}' \mathbf{A} \boldsymbol{\mu} / \sigma^2$ .

$$\text{Theorem: } X \sim N(0, 1) \\ u \sim \chi^2(m) \\ X \text{ and } u \text{ independent} \quad \left\{ \begin{array}{l} \frac{X}{\sqrt{u/m}} \sim \chi^2(m) \\ \uparrow \\ \text{degree of freedom} \end{array} \right.$$

3. Let  $B$  be a  $q \times k$  matrix, and let  $W$  be the linear form given by  $W = By$ . The quadratic form  $U = \underbrace{y' A y}_{\text{MSRej}}$  and  $W$  are independent if

$$BVA = 0.$$

$$\hat{\beta}_1 = \sum c_i y_i \\ (\text{function of } y)$$

Note: If  $V = \sigma^2 I$ , then  $U$  and  $W$  are independent if  $BA = 0$ .

4. Let  $B$  be a  $k \times k$  matrix. Let  $V = y' B y$ . The two quadratic forms,  $U$  and  $V$ , are independent if

$$AVB = 0.$$

Note: If  $V = \sigma^2 I$ , then  $U$  and  $V$  are independent if  $AB = 0$ .

$$u \sim \chi^2(n)$$

$$v \sim \chi^2(n)$$

$$\frac{y_m}{v/n} \sim F_{m, n}$$

The End