# Chapter 2  Categorical variables

Model $\underset{\sim}{Y} = \underset{\sim}{X}\underset{\sim}{\beta} + \underset{\sim}{e}$    $\underset{\sim}{e} \overset{\text{Assumption}}{\sim} N(0, \sigma^2 \underset{\sim}{I})$

$\hat{\underset{\sim}{\beta}} = (\underset{\sim}{X}^T\underset{\sim}{X})^{-1}\underset{\sim}{X}^T\underset{\sim}{Y}$    $\hat{\underset{\sim}{\beta}} \sim N(\underset{\sim}{\beta}, \sigma^2 (\underset{\sim}{X}^T\underset{\sim}{X})^{-1})$

$\hat{\sigma}^2 = \dfrac{\text{Res S.S.}}{n - p'}$    $\dfrac{\text{Res S.S.}}{\sigma^2} \sim \chi^2(n-p')$

where Res S.S. $= \underset{\sim}{Y}^T(\underset{\sim}{I} - \underset{\sim}{X}(\underset{\sim}{X}^T\underset{\sim}{X})^{-1}\underset{\sim}{X}^T)\underset{\sim}{Y}$

Chapter 1

$\Rightarrow$ C.I. of $\beta$ , $t$-test , $F$-test

## One categorical variable (One-way ANOVA)
$\uparrow$

e.g. p.4

Fit $y$ on group $= 1, 2, 3, 4, 5$

Model $y = \beta_0 + \beta_1 x + e$    $\longleftarrow$ Lack-of-fit

| group | Model | E(y) |
|-------|-------|------|
| 1 | $y = \beta_0 + \beta_1 + e$ | $\beta_0 + \beta_1 = \mu_1$ |
| 2 | $y = \beta_0 + 2\beta_1 + e$ | $\beta_0 + 2\beta_1 = \mu_2$ |
| 3 | $y = \beta_0 + 3\beta_1 + e$ | $\beta_0 + 3\beta_1 = \mu_3$ |
| 4 | $y = \beta_0 + 4\beta_1 + e$ | $\beta_0 + 4\beta_1 = \mu_4$ |
| 5 | $y = \beta_0 + 5\beta_1 + e$ | $\beta_0 + 5\beta_1 = \mu_5$ |

$\mu_2 - \mu_1 = \mu_3 - \mu_2 = \mu_4 - \mu_3 =$

$\mu_5 - \mu_4$

$\Downarrow$

It is not a general model

## general model    Dummy variable (indictor variable)

| group | $g_1$ | $g_2$ | $g_3$ | $g_4$ | $g_5$ |
|-------|-------|-------|-------|-------|-------|
| 1 | 1 | 0 | 0 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 0 | 0 |
| 4 | 0 | 0 | 0 | 1 | 0 |
| 5 | 0 | 0 | 0 | 0 | 1 |

group $= i$ , $\Rightarrow$ $g_i = 1$    $i = 1, 2, 3, 4, 5$

$$\underset{n\times 6}{\underset{\sim}{X}} = \begin{pmatrix} & \beta_0 & \alpha_1 & \alpha_2 & \alpha_3 & \alpha_4 & \alpha_5 \\ & \vdots & 1 & 0 & 0 & 0 & 0 \\ & & & & & & \\ & & 0 & 1 & 0 & 0 & 0 \\ & & 0 & 0 & 1 & 0 & 0 \\ & & 0 & 0 & 0 & 1 & 0 \\ & & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

group = 1

group 2

group = 3

group = 4

group = 5

# of levels of group = $\overset{m}{\cancel{5}}$

$\Rightarrow$ create $(\overset{m}{\cancel{5}}-1)$ dummy variables

$\det\left(\underset{\sim}{X}^T\underset{\sim}{X}\right) = 0$ $\qquad$ $\underset{\sim}{X}$ is singular matrix

we need to add one constraint. e.g. delete the last column

## Model I (Regression model)

$+ \cdots + \alpha_{m-1}\, g_{i,m-1}$

$$y_i = \beta_0 + \alpha_1 g_{i1} + \alpha_2 g_{i2} + \alpha_3 g_{i3} + \alpha_4 g_{i4} + e_i \qquad i = 1, \longrightarrow, n$$

$\alpha_k$ – regression coeff. of $g_k$ $\qquad k = 1, 2, \cdots, m-1$

Analysis of variance

(ANOVA) model

$$y_{ij} = \beta_0 + \alpha_i + e_{ij} \qquad \begin{array}{l} i = 1, \longrightarrow, m \\ j = 1, \longrightarrow, n_i \end{array}$$

| Group | Model | | |
|---|---|---|---|
| 1 | $E(y_{1j}) = \beta_0 + \alpha_1 + \cancel{e_{1j}}$ | $\mu_1$ | $j = 1, \cdots, n_1$ |
| 2 | $E(y_{2j}) = \beta_0 + \alpha_2 + \cancel{e_{2j}}$ | $\mu_2$ | $j = 1, \longrightarrow, n_2$ |
| 3 | $E(y_{3j}) = \beta_0 + \alpha_3 + \cancel{e_{3j}}$ | $\mu_3$ | $j = 1, \longrightarrow, n_3$ |
| 4 | $E(y_{4j}) = \beta_0 + \alpha_4 + \cancel{e_{4j}}$ | $\mu_4$ | $j = 1, \longrightarrow, n_4$ |
| 5 | $E(y_{5j}) = \beta_0 + \cancel{\alpha_5} + \cancel{e_{5j}}$ | $\mu_5$ | $j = 1, \longrightarrow, n_5$ |

$$\Rightarrow \qquad y_{ij} = \mu_i + e_{ij} \qquad \begin{array}{l} i = 1, \longrightarrow, m \\ j = 1, \longrightarrow, n_i \end{array}$$

$$\beta_0 + \alpha_1 = \mu_1 \qquad\qquad \alpha_1 = \mu_1 - \mu_5$$
$$\beta_0 + \alpha_2 = \mu_2 \qquad\qquad \alpha_2 = \mu_2 - \mu_5$$
$$\beta_0 + \alpha_3 = \mu_3 \quad\Longrightarrow\quad \alpha_3 = \mu_3 - \mu_5$$
$$\beta_0 + \alpha_4 = \mu_4 \qquad\qquad \alpha_4 = \mu_4 - \boxed{\mu_5} \leftarrow \text{group} = 5 \text{ is a}$$
$$\beta_0 \quad\;\; = \mu_5 \qquad\qquad\qquad\qquad\qquad \text{reference group}$$

$\leftarrow$ Model = — boxed: model without intercept

$$y_{ij} = \mu_i + \ell_{ij}$$
$$i = 1, \dots, 5/m$$
$$j = 1, \dots, n_i / 6$$

$$X_{n \times 5} =
\begin{pmatrix}
\mu_1 & \mu_2 & \mu_3 & \mu_4 & \mu_5 \\
\vdots & 0 & 0 & 0 & 0 \\
\hline
0 & 0 & 0 & 0 & 0 \\
\hline
0 & 0 & 0 & 0 & 0 \\
\hline
0 & 0 & 0 & 0 & 0 \\
\hline
0 & 0 & 0 & 0 & 
\end{pmatrix}
\begin{matrix}
\text{group}=1 \\ \text{group}=2 \\ \text{group}=3 \\ \text{group}=4 \\ \text{group}=5
\end{matrix}$$

$$Y = \begin{pmatrix} y_{11} \\ y_{1n_1} \\ y_{21} \\ y_{2n_2} \\ \vdots \\ y_{51} \\ y_{5n_5} \end{pmatrix}
\begin{matrix} \text{group}=1 \\ \text{group}=2 \\ \text{group}=3 \\ \text{group}=4 \\ \text{group}=5 \end{matrix}
= \begin{pmatrix} 551 \\ 632 \\ 595 \\ 517 \\ \vdots \\ 563 \\ 679 \end{pmatrix}
\begin{matrix} \text{group}=1 \\ \text{group}=2 \\ \text{group}=5 \end{matrix}$$

$$X^T X_{5 \times 5} =
\begin{pmatrix}
n_1 & 0 & 0 & 0 & 0 \\
0 & n_2 & 0 & 0 & 0 \\
0 & 0 & n_3 & 0 & 0 \\
0 & 0 & 0 & n_4 & 0 \\
0 & 0 & 0 & 0 & n_5
\end{pmatrix}$$

$$(X^T X)^{-1} =
\begin{pmatrix}
\frac{1}{n_1} & 0 & 0 & 0 & 0 \\
0 & \frac{1}{n_2} & 0 & 0 & 0 \\
0 & 0 & \frac{1}{n_3} & 0 & 0 \\
0 & 0 & 0 & \frac{1}{n_4} & 0 \\
0 & 0 & 0 & 0 & \frac{1}{n_5}
\end{pmatrix}$$

$$X^T Y =
\begin{pmatrix}
\sum_{j=1}^{n_1} y_{1j} \\
\sum_{j=1}^{n_2} y_{2j} \\
\vdots \\
\sum_{j=1}^{n_5} y_{5j}
\end{pmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y
= \begin{pmatrix} \bar{y}_{1\cdot} \\ \bar{y}_{2\cdot} \\ \bar{y}_{3\cdot} \\ \bar{y}_{4\cdot} \\ \bar{y}_{5\cdot} \end{pmatrix}$$

$$\hat{\beta} = \begin{pmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \\ \hat{\mu}_3 \\ \hat{\mu}_4 \\ \hat{\mu}_5 \end{pmatrix}$$

$$\hat{\mu}_i = \bar{y}_{i\cdot} \quad i = 1, \dots, m$$

③

## Regression model

$$y_i = \beta_0 + \alpha_1 g_{i1} + \alpha_2 g_{i2} + \alpha_3 g_{i3} + \alpha_4 g_{i4} + e_i$$

## ANOVA model

$$y_{ij} = \mu_i + e_{ij} \qquad \begin{array}{l} i = 1, \cdots, 5 \\ j = 1, \cdots, n_i \end{array}$$

$$\boxed{\begin{aligned}
\hat{\beta}_0 + \hat{\alpha}_1 &= \hat{\mu}_1 = \bar{y}_{1\cdot} \\
\hat{\beta}_0 + \hat{\alpha}_2 &= \hat{\mu}_2 = \bar{y}_{2\cdot} \\
\hat{\beta}_0 + \hat{\alpha}_3 &= \hat{\mu}_3 = \bar{y}_{3\cdot} \\
\hat{\beta}_0 + \hat{\alpha}_4 &= \hat{\mu}_4 = \bar{y}_{4\cdot} \\
\hat{\beta}_0 &= \hat{\mu}_5 = \bar{y}_{5\cdot}
\end{aligned}}$$

$$\Rightarrow \quad \hat{\beta}_0 = \bar{y}_{5\cdot}, \quad \hat{\alpha}_1 = \bar{y}_{5\cdot} - \bar{y}_{1\cdot}, \quad \cdots$$

---

$$\text{Res S.S.} = \underset{\sim}{Y}^T \underset{\sim}{Y} - \hat{\underset{\sim}{\beta}}^T X^T \underset{\sim}{Y}$$

$$= \sum_{i=1}^{m} \boxed{\sum_{j=1}^{n_i} y_{ij}^2} - (\bar{y}_1, \bar{y}_2, \cdots, \bar{y}_m) \boxed{\begin{pmatrix} \sum_{j=1}^{n_1} y_{1j} \\ \vdots \\ \sum_{j=1}^{n_m} y_{mj} \end{pmatrix}}$$

For $i = 1$

$$\sum_{j=1}^{n_1} y_{1j}^2 - \bar{y}_{1\cdot} \underbrace{\left( \sum_{j=1}^{n_1} y_{1j} \right)}_{n_1 \bar{y}_1} \leftarrow n_1 \bar{y}_1$$

$$= \sum_{j=1}^{n} y_{1j}^2 - n \bar{y}_{1\cdot}^2$$

$$= \sum_{j=1}^{n} (y_{1j} - \bar{y}_{1\cdot})^2 \qquad \text{Pure Error S.S.}$$

$$\Rightarrow \text{Res S.S.} = \boxed{\sum_{i=1}^{m} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}$$

MATH 3423

$$\frac{\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}{\sigma^2} \sim \chi^2_{(n_i - 1)}$$

$$\frac{\text{Res S.S.}}{\sigma^2} \sim \chi^2_{\underbrace{\sum_{i=1}^{m}(n_i - 1)}_{n - m}} \quad \Rightarrow \quad \hat{\sigma}^2 = \frac{\text{Res S.S.}}{n - m} = \frac{\sum_{i=1}^{m} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}{n - m}$$

④

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

$$= \begin{pmatrix} \frac{1}{n_1} & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{n_2} & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{n_3} & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{n_4} & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{n_5} \end{pmatrix} \sigma^2$$

$$\beta = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \\ \mu_5 \end{pmatrix}$$

$$\text{Var}(\hat{\mu}_i) = \frac{\sigma^2}{n_i}$$

$$\Rightarrow \text{Var}(\bar{y}_i) = \frac{\sigma^2}{n_i}$$

C.I. for $\mu_i$

$$\bar{y}_{i.} \pm t_{\alpha/2,\, n-m} \,\hat{\sigma} \sqrt{\frac{1}{n_i}}$$

$$H_0 = \mu_i = \mu_{i_0}$$

$$t = \frac{\bar{y}_{i.} - \mu_{i_0}}{\hat{\sigma}/\sqrt{n_i}}$$

---

Fitted model : $y$ on group $= 1, 2, 3, 4, 5$

$$y_i = \beta_0 + \beta_1 * \text{group } i + e \qquad i = 1, \dots, n$$

$$1, 2, 3, 4, 5$$

Section 6 of Chapter 1

$H_0$ : No lack of fit

$H_1 = $

$$F = \frac{\boxed{\text{lack of fit S.S.}} \,/\, \boxed{\text{d.f} \mid \text{lack of fit S.S.}}}{\boxed{\text{Pure Error S.S}} \,/\, \boxed{\text{d.f} \mid \text{pure error s.s.}}}$$

Res S.S. | fitted model − Pure Error S.S.

$n - 2$

− d.f. of Res S.S. | fitted model

− d.f. of Pure Error S.S.

$n - m$

$\boxed{\text{Pure Error S.S}} = $
$$\overset{n}{\underset{i=1}{\sum}} \overset{n_i}{\underset{j=1}{\sum}} (y_{ij} - \bar{y}_i)^2$$
$$= \text{Res S.S.} \mid_{\text{ANOVA}}$$

$\boxed{\text{d.f} \mid \text{pure error s.s.}} = $
$$n - m$$

$$H_0 = \mu_1 = \mu_2 = \cdots = \mu_m = \mu$$

$$\text{vs } H_1 : \text{at least } \underset{two}{\sout{one}} \text{ of } \mu\text{'s are not equal}$$

ANOVA model

$$\Updownarrow \quad y_{ij} = \mu_i + e_{ij}$$

$$\text{Under } H_0, \quad y_{ij} = \mu + e_{ij}$$

$$H_0 : \alpha_1 = \alpha_2 = \cdots = \alpha_{m-1} = 0$$

Regression model

$$y_i = \beta_0 + \alpha_1 g_{i1} + \cdots + \alpha_{m-1} g_{i,m-1} + e_i$$

$$\text{Under } H_0, \quad y_i = \beta_0 + e_i$$

Section 4 of Chapter 1

Total S.S. = Reg S.S. + Res S.S.

$$\sum_{i=1}^{m} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^{m} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.} + \bar{y}_{i.} - \bar{y}_{..})^2$$

overall means of $y_{ij}$

$$= \sum_{i=1}^{m} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.})^2 + \sum_{i=1}^{m} \sum_{j=1}^{n_i} (\bar{y}_{i.} - \bar{y}_{..})^2$$

$$\Vert$$

$$\sum_{i=1}^{m} n_i (\bar{y}_{i.} - \bar{y}_{..})^2$$

$$\frac{\text{Total S.S}}{\sigma^2} \sim \chi^2_{(n-1, \lambda)}$$

$$\frac{\text{Res S.S.}}{\sigma^2} \sim \chi^2_{(n-m)}$$

$$\text{Reg S.S.} \sim \chi^2_{(m-1, \lambda)}$$

$$\Rightarrow \quad F = \frac{\text{Reg S.S.}/(m-1)}{\text{Res S.S.}/(n-m)} \overset{\text{under } H_0}{\sim} F_{(m-1, n-m)}$$

$$\text{Reg S.S.} = \sum_{i=1}^{m} n_i (\bar{y}_{i.} - \bar{y}_{..})^2 \qquad N = \sum_{i=1}^{m} n_i$$

$$= \text{total \# of observations}$$

$$= \sum_{i=1}^{m} n_i \left( \frac{T_{i.}}{n_i} - \frac{T_{..}}{N} \right)^2$$

$$= \sum_{i=1}^{m} n_i \left( \frac{T_{i.}^2}{n_i^2} + \frac{T_{..}^2}{N^2} - 2 \frac{T_{i.}}{n_i} \frac{T_{..}}{N} \right) \qquad \Vert T_{..}$$

$$= \sum_{i=1}^{m} \frac{T_{i.}^2}{n_i} + \frac{T_{..}^2}{N^2} \underbrace{\sum_{i=1}^{m} n_i}_{n} - 2 \frac{T_{..}}{N} \boxed{\sum_{i=1}^{m} n_i \cdot \frac{T_{i.}}{n_i}}$$

$$= \sum_{i=1}^{m} \frac{T_{i.}^2}{n_i} - \frac{T_{..}^2}{N}$$

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_m = \mu \iff H_0 : \alpha_1 = \alpha_2 = \cdots = \alpha_{m-1} = 0$$

| Source of Variation of | Sum of Squares | Degrees of freedom | Mean Square | Computed $f$ |
|---|---|---|---|---|
| Model | $\sum_{i=1}^{m} n_i(\bar{y}_{i.} - \bar{y}_{..})^2$ | $m - 1$ | $\dfrac{\sum_{i=1}^{m} n_i(\bar{y}_{i.} - \bar{y}_{..})^2}{m-1}$ | $\dfrac{(\sum_{i=1}^{m} n_i - m)\sum_{i=1}^{m} n_i(\bar{y}_{i.} - \bar{y}_{..})^2}{(m-1)\sum_{i=1}^{m}\sum_{j=1}^{n_i}(y_{ij} - \bar{y}_{i.})^2}$ |
| Error | $\sum_{i=1}^{m}\sum_{j=1}^{n_i}(y_{ij} - \bar{y}_{i.})^2$ | $\left(\sum_{i=1}^{m} n_i\right) - m$ $\quad \overset{n}{=}$ | $\dfrac{\sum_{i=1}^{m}\sum_{j=1}^{n_i}(y_{ij} - \bar{y}_{i.})^2}{\sum_{i=1}^{m} n_i - m}$ | |
| Total | $\sum_{i=1}^{m}\sum_{j=1}^{n_i}(y_{ij} - \bar{y}_{..})^2$ | $\sum_{i=1}^{m} n_i - 1$ | | |

The advantages of choosing equal sample sizes over the choice of unequal sample sizes are: 1) the $f$ ratio is insensitive to slight departures from the assumption of equal variances for the $m$ populations when the sample are of equal sizes; and 2) the choice of equal sample size minimizes the probability of committing a type II error.

**Example**

*$S_i^2$ – sample variance of $y_{ij}$ in the $i$th group*

*$X_i$ categorical variable*

*$S_T^2$ — Sample variance of all $y_{ij}$*

*repeated measurements*

| | Group | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | |
| $y_{11}$ | 551 | 595 | 639 | 417 | 563 | $y_{51}$ |
| $y_{12}$ | 457 | 580 | 615 | 449 | 631 | |
| | 450 | 508 | 511 | 517 | 522 | |
| | 731 | 583 | 573 | 438 | 613 | |
| | 499 | 633 | 648 | 415 | 656 | $y_{56}$ |
| $y_{16}$ | 632 | 517 | 677 | 555 | 679 | |
| Total | 3320 | 3416 | 3663 | 2791 | 3664 | 16854 |
| Mean | 553.33 | 569.33 | 610.50 | 465.19 | 610.67 | 561.80 |

*$S_{yy}^{(i)}$ = Sum of squares of $y_{ij}$ in the $i$th group*

**Solution**

$S_1^2 = 12,133.8667, \quad S_2^2 = 2,302.6667, \quad S_3^2 = 3593.5, \quad S_4^2 = 3,318.5667, \quad S_5^2 = 3,455.4667,$
$S_T^2 = 7,219.8897$

$$\sum_{i=1}^{m}\sum_{j=1}^{n_i}(y_{ij} - \bar{y}_{i.})^2 = \sum_{i=1}^{m}(n_i - 1)S_i^2$$

| Source of Variation of | Sum of Squares | Degrees of freedom | Mean Square | Computed $f$ |
|---|---|---|---|---|
| Group | 85,356 | 4 | 21,339 | 4.30 |
| Error | 124,021 | 25 $= \sum_{i=1}^{m}(n_i - 1) = n - m$ | 4,961 | |
| Total | 209,377 | 29 | | |

① Reg S.S. = Total S.S. − Res S.S.

② $\sum_{i=1}^{m} \dfrac{T_{i.}^2}{n_i} - \dfrac{T_{..}^2}{n}$

$\dfrac{21339}{4961} = 4.30$

$C.V. = F_{0.05, 4, 2} = 2.76$

$30 - 1$

The critical value $f_{0.05}(4, 25) = 2.76$. Thus, $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$ is rejected.

$$\sum_{i=1}^{m}\sum_{j=1}^{n_i}(y_{ij} - \bar{y}_{..})^2 = (n-1)S_T^2 = S_{yy}$$

*overall sample size*

$\Rightarrow$ Reject $H_0$

What is the reason of rejecting the null hypothesis?

⑰