Indicator (Dummy) Variables – Part III

Johns Hopkins Engineering

625.461 Statistical Models and Regression

Module 6 – Lecture 6B



Regression Approach to Analysis of Variance (ANOVA)

ANOVA is a technique frequently used to analyze data from planned or designed experiments.

$$y_{ij} = \mu + \tau_i + \varepsilon_{ij}, \quad i = 1, 2, ..., k, \quad j = 1, 2, ..., n$$

 y_{ij} : jth obs for ith treatment (or factor level)

 μ : grand mean

 τ_i : the effect of *i*th treatment

 ε_{ij} is NIND(0, σ^2)

In the balanced case (i.e., equal number of observations per treatment), $\sum_{i=1}^{k} \tau_i = 0$

The mean of the *i*th treatment is

$$\mu_i = \mu + \tau_i$$
 , $i = 1, ..., k$

The average of *n* observations in the *i*th treatment:

$$\overline{y}_{i.} = \frac{1}{n} \sum_{j=1}^{n} y_{ij}, \quad i = 1, 2, ..., k$$

The grand average is

$$\bar{y}_{..} = \frac{1}{kn} \sum_{i=1}^{k} \sum_{j=1}^{n} y_{ij}$$

How to use indicator variables to perform regression analysis to obtain ANOVA?

Without loss of generality, let us assume that k = 3. Define

$$x_1 = \begin{cases} 1 & \text{if the observation is from treatment 1} \\ 0 & \text{otherwise} \end{cases}$$

$$x_2 = \begin{cases} 1 & \text{if the observation is from treatment 2} \\ 0 & \text{otherwise} \end{cases}$$

1-Way ANOVA vs. Regression Model

$$y_{ij} = \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2j} + \varepsilon_{ij}, \quad i = 1, 2, 3, \quad j = 1, 2, ..., n$$

For treatment 1:

$$y_{1j} = \beta_0 + \beta_1(1) + \beta_2(0) + \varepsilon_{1j} = \beta_0 + \beta_1 + \varepsilon_{1j}$$

$$\beta_0 + \beta_1 = \mu_1$$

1-Way ANOVA vs. Regression Model

$$y_{ij} = \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2j} + \varepsilon_{ij}, \quad i = 1, 2, 3, \quad j = 1, 2, ..., n$$

For treatment 2:

$$y_{2j} = \beta_0 + \beta_1(0) + \beta_2(1) + \varepsilon_{2j} = \beta_0 + \beta_2 + \varepsilon_{2j}$$
$$\beta_0 + \beta_2 = \mu_2$$

For treatment 3:

$$y_{3j} = \beta_0 + \beta_1(0) + \beta_2(0) + \varepsilon_{3j} = \beta_0 + \varepsilon_{3j}$$

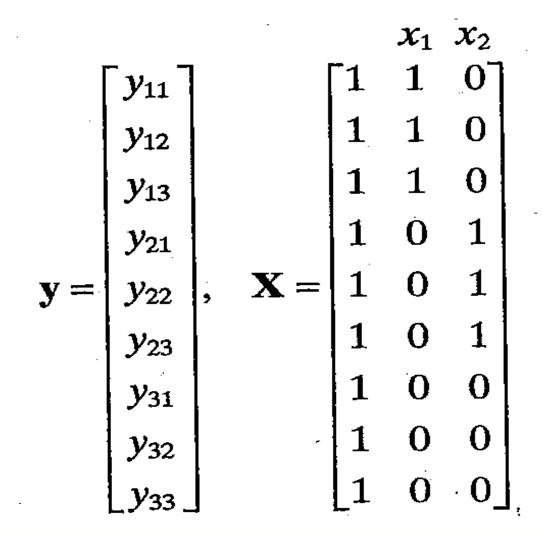
 $\beta_0 = \mu_3$

1-Way ANOVA vs. Regression Model

Consequently,

$$eta_0 = \mu_3 \ eta_1 = \mu_1 - \mu_3 \ eta_2 = \mu_2 - \mu_3$$

Assume that n = 3



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

$$\begin{bmatrix} 9 & 3 & 3 \\ 3 & 3 & 0 \\ 3 & 0 & 3 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} y_{..} \\ y_{1.} \\ y_{2.} \end{bmatrix}$$

$$\hat{\boldsymbol{\beta}}_0 = \overline{y}_{..} - \overline{y}_{1.} - \overline{y}_{2.} = \overline{y}_{3.}, \quad \hat{\boldsymbol{\beta}}_1 = \overline{y}_{1.} - \overline{y}_{3.}, \quad \hat{\boldsymbol{\beta}}_2 = \overline{y}_{2.} - \overline{y}_{3.}$$

$$SS_{R}(\hat{\beta}_{0}, \hat{\beta}_{1}, \hat{\beta}_{2}) = \hat{\beta}' \mathbf{X}' \mathbf{y} = [\overline{y}_{3}, \overline{y}_{1}, -\overline{y}_{3}, \overline{y}_{2}, -\overline{y}_{3},]\begin{bmatrix} y_{.} \\ y_{1} \\ y_{2} \end{bmatrix}$$

$$= y_{.} \overline{y}_{3} + y_{1} (\overline{y}_{1}, -\overline{y}_{3}) + y_{2} (\overline{y}_{2}, -\overline{y}_{3})$$

$$= (y_{1} + y_{2} + y_{3}) \overline{y}_{3} + y_{1} (\overline{y}_{1}, -\overline{y}_{3}) + y_{2} (\overline{y}_{2}, -\overline{y}_{3})$$

$$= \overline{y}_{1} y_{1} + \overline{y}_{2} y_{2} + \overline{y}_{3} y_{3}$$

$$= \sum_{i=1}^{3} \frac{y_{i}^{2}}{3}$$

$$SS_{Res} = \sum_{i=1}^{3} \sum_{j=1}^{3} y_{ij}^{2} - SS_{R}(\beta_{0}, \beta_{1}, \beta_{2})$$

$$= \sum_{i=1}^{3} \sum_{j=1}^{3} y_{ij}^{2} - \sum_{i=1}^{3} \frac{y_{i.}^{2}}{3}$$

$$= \sum_{i=1}^{3} \sum_{j=1}^{3} (y_{ij} - \overline{y}_{i.})^{2}$$

To test H_0 : $\tau_1 = \tau_2 = \tau_3 = 0$, equivalently,

$$\beta_0 = \mu$$
, $\beta_1 = 0$, $\beta_2 = 0$, the reduced model: $y_{ij} = \beta_0 + \varepsilon_{ij}$

$$SS_{R}(\beta_{1}, \beta_{2} | \beta_{0}) = SS_{R}(\beta_{0}, \beta_{1}, \beta_{2}) - SS_{R}(\beta_{0})$$

$$= \sum_{i=1}^{3} \frac{y_{i.}^{2}}{3} - \frac{y_{..}^{2}}{9}$$

$$= 3\sum_{i=1}^{3} (\overline{y}_{i.} - \overline{y}_{..})^{2}$$

$$F_{0} = \frac{SS_{R}(\beta_{1}, \beta_{2}|\beta_{0})/2}{SS_{Res}/6}$$

$$= \frac{3\sum_{i=1}^{3} (\overline{y}_{i.} - \overline{y}_{..})^{2}/2}{\sum_{i=1}^{3} \sum_{j=1}^{3} (y_{ij} - \overline{y}_{i.})^{2}/6}$$

$$= \frac{MS_{Treatments}}{MS_{Dec}}$$

TABLE 8.4	One-Way Analysis of Variance			
Degrees of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F_0
Treatments	$n\sum_{i=1}^k (\overline{y}_{i.} - \overline{y}_{})^2$	<i>k</i> – 1	$\frac{SS_{\text{Treatments}}}{k-1}$	$\frac{MS_{\mathrm{Treatments}}}{MS_{\mathrm{Res}}}$
Error	$\sum_{i=1}^{k} \sum_{j=1}^{n} (y_{ij} - \overline{y}_{i.})^{2}$	k(n-1)	$\frac{SS_{\mathrm{Res}}}{k(n-1)}$	
Total	$\sum_{i=1}^{k} \sum_{j=1}^{n} (y_{ij} - \bar{y}_{})^{2}$	<i>kn</i> − 1		·

The regression approach using appropriate dummy variables is identical to 1-way ANOVA

