

XIII. Analysis of Variance

10

- only categorical covariates

$$Z = (z_1, \dots, z_p)^T, \quad z_j \in \mathcal{Z}_j = \text{finite, low cardinality}$$

Main usage agricultural

- designed (industrial) experiments
- controlled clinical trials

$Z \equiv$ p factors that may influence outcome

e.g. $z_1 =$ type of production machine (A, B, C)

$z_2 =$ production temperature
(experiment done at prespecified temp. values of 20°C, 25°C, 30°C)

$z_3 =$ used amount of some ingredient
(few prespecified values used at experiment)

→ covariates are mostly fixed, not random, but this does not matter for theory (we condition by covariate values anyway)

13.1 One-way classification

1

≡ REPETITION from Mathematical Statistics 1 and from Chapter 4

$$p=1, Z \in \mathcal{Z} = \{1, \dots, G\}$$

2

$$E(Y|Z=g) = m(g) =: mg, \quad g=1, \dots, G$$

without loss of generality

→ data sorted by Z values

→ double subscript used

$$Z = \begin{pmatrix} z_1 \\ \vdots \\ z_{m_1} \\ z_{m_1+1} \\ \vdots \\ z_{m_1+m_2} \\ \vdots \\ z_{m_1+m_2+\dots+m_{G-1}+1} \\ \vdots \\ z_{m_1+\dots+m_{G-1}+m_G} \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ 1 \\ 2 \\ \vdots \\ 2 \\ \vdots \\ \vdots \\ G \\ \vdots \\ G \end{pmatrix}, \quad Y = \begin{pmatrix} y_{11} \\ \vdots \\ y_{1,m_1} \\ y_{2,1} \\ \vdots \\ y_{2,m_2} \\ \vdots \\ \vdots \\ y_{G,1} \\ \vdots \\ y_{G,m_G} \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_2 \\ \vdots \\ \vdots \\ y_g \end{pmatrix}$$

→ by linear model, we use "model"

$$E(Y|Z) = \begin{pmatrix} m_1 \mathbb{1}_{m_1} \\ \vdots \\ m_G \mathbb{1}_{m_G} \end{pmatrix}$$

3

ASSUMPTION (if LM used):

$$\text{var}(Y|Z) = \sigma^2 \cdot I_n$$

(homoscedasticity)

We will also assume:

$$P(m_g > 0) = 1 \quad \forall g=1, \dots, G$$

2

13.1.1 Parameters of interest

4

(1) Differences between the group means

$$\theta_{g,h} = m_g - m_h, \quad g, h = 1, \dots, G, \quad g \neq h$$

↑ ↑
elements of a vector $E(Y|Z)$

$$\rightarrow \text{LSE of } \theta_{g,h}, \quad \hat{\theta}_{g,h} = \hat{Y}_g - \hat{Y}_h$$

↑ ↑
fitted values for obs.
with $Z=g$ and $Z=h$

→ principal null hypothesis
to be tested

$$H_0: m_1 = \dots = m_G$$

$$\theta_{g,h} = 0, \quad g, h = 1, \dots, G, \quad g \neq h$$

(2) Factor effects

5

Def 13.1 Factor effects in a one-way classification

By factor effects in case of a one-way classification we understand the quantities defined as

$$\eta_g = m_g - \bar{m}, \quad g = 1, \dots, G, \quad \text{where}$$
$$\bar{m} = \frac{1}{G} \sum_{h=1}^G m_h \quad \text{is the mean of the group means.}$$

NOTE: $\bar{m} = \frac{1}{G} \sum_{h=1}^G E(Y|Z=h) = EY$

↑ iff $P(Z=h) = \frac{1}{G} \quad \forall h=1, \dots, G$

\bar{m} = mean result of experiment if each setting is used equally often

3

13.1.2 One-way ANOVA model

6

Regression space

$$\left\{ \begin{pmatrix} m_1 \\ \vdots \\ m_G \end{pmatrix} : m_1, \dots, m_G \in \mathbb{R} \right\} \subseteq \mathbb{R}^n$$

$$\text{if } m_g > 0 \forall g \Rightarrow \text{rec-dim} = G$$

Full-rank parameterizations

7

$$m_g = \beta_0 + C_g^T \beta^z, \quad g = 1, \dots, G$$

$$\beta = (\beta_0, (\beta^z)^T)^T$$
$$(\beta_1^z, \dots, \beta_{G-1}^z)^T$$

$$C = \begin{pmatrix} C_1^T \\ \vdots \\ C_G^T \end{pmatrix} \equiv \text{chosen (pseudo)contrast matrix with } G-1 \text{ columns}$$

NOTE: With $C = \begin{pmatrix} 1 & \dots & 0 \\ \vdots & & \\ 0 & \dots & 1 \\ -1 & \dots & -1 \end{pmatrix} \equiv \text{sum contrasts}$

$$\alpha_0 := \beta_0 = \frac{1}{G} \sum_{g=1}^G m_g = \bar{m}$$

factor effects

$$\left\{ \begin{aligned} \alpha_g &:= \beta_g = m_g - \bar{m} = \eta_g, \quad g = 1, \dots, G-1 \\ \alpha_G &:= -\sum_{g=1}^{G-1} \beta_g = m_G - \bar{m} = \eta_G \end{aligned} \right.$$

4

13.1.3 Least squares estimation

8

Lemma 13.1 Least squares estimation in one-way ANOVA

short:

$$\hat{m}_g = \hat{y}_{g \cdot j} = \frac{1}{n_g} \sum_{l=1}^{n_g} y_{gl} =: \bar{y}_g$$

$g=1, \dots, G, j=1, \dots, n_g$

$$\hat{m} = \begin{pmatrix} \hat{m}_1 \\ \vdots \\ \hat{m}_G \end{pmatrix} = \begin{pmatrix} \bar{y}_{1 \cdot} \\ \vdots \\ \bar{y}_{G \cdot} \end{pmatrix} \quad \hat{y} = \begin{pmatrix} \bar{y}_{1 \cdot} \mathbb{1}_{n_1} \\ \vdots \\ \bar{y}_{G \cdot} \mathbb{1}_{n_G} \end{pmatrix}$$

Under normality

$$\hat{m} | Z \sim N_G(m, \sigma^2 V), \quad V = \begin{pmatrix} \frac{1}{n_1} & 0 \\ 0 & \frac{1}{n_G} \end{pmatrix}$$

Proof: Possible full-rank parameterization is

$$E(Y|Z) = \begin{pmatrix} m_1 \mathbb{1} \\ \vdots \\ m_G \mathbb{1} \end{pmatrix} = \begin{pmatrix} \mathbb{1}_{n_1} & 0 \\ 0 & \mathbb{1}_{n_G} \end{pmatrix} \cdot \underbrace{\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_G \end{pmatrix}}_{\beta}$$

i.e. $\beta_g = m_g, g=1, \dots, G$

$$\Rightarrow (X^T X) = \begin{pmatrix} n_1 & 0 \\ 0 & n_G \end{pmatrix}, \quad (X^T X)^{-1} = \begin{pmatrix} \frac{1}{n_1} & 0 \\ 0 & \frac{1}{n_G} \end{pmatrix}$$

$$X^T Y = \begin{pmatrix} \sum_j y_{1j} \\ \vdots \\ \sum_j y_{Gj} \end{pmatrix}$$

$$\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T Y = \begin{pmatrix} \bar{y}_{1 \cdot} \\ \vdots \\ \bar{y}_{G \cdot} \end{pmatrix} = \hat{m}$$

$$\hat{Y} = X\hat{\beta} = \begin{pmatrix} \bar{y}_1 \\ \vdots \\ \bar{y}_g \\ \vdots \\ \bar{y}_g \\ \vdots \\ \bar{y}_g \end{pmatrix}$$

□

$$\text{var}(\hat{m} | Z) = \text{var}(\hat{\beta} | Z) = \sigma^2 (X^T X)^{-1} = \sigma^2 \begin{pmatrix} \frac{1}{n_1} & 0 \\ 0 & \frac{1}{n_g} \end{pmatrix}$$

↓

normality \equiv general LSE theory

□

LSE of regression coefficients and their linear combinations

Full-rank parameterization:

$$m_g = \beta_0 + C_g^T \beta^z \quad \rightarrow \quad m = \beta_0 \mathbb{1} + C \beta^z$$

$$\begin{pmatrix} \bar{y}_{10} \\ \vdots \\ \bar{y}_{G0} \end{pmatrix} = \hat{m} = (\text{subvector of fitted values}) = \hat{\beta}_0 \mathbb{1} + C \hat{\beta}^z$$
$$\hat{y} = X \hat{\beta}$$

$$\hat{m} = \underbrace{(1, C)}_{\text{invertible}} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}^z \end{pmatrix}$$

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}^z \end{pmatrix} = (1, C)^{-1} \hat{m}$$

↑
vector of group sample means

- in any parameterization

• $\hat{\beta}$ = linear combination of group sample means

• the same holds for any linear combination of β 's, e.g.

$$* \quad \theta_{gn} = m_g - m_n \quad \rightarrow \quad \hat{\theta}_{gn} = \hat{m}_g - \hat{m}_n = \bar{y}_{g0} - \bar{y}_{n0}$$

$$* \quad \eta_g = m_g - \bar{m} \quad \rightarrow \quad \hat{\eta}_g = \hat{m}_g - \frac{1}{G} \sum_{h=1}^G \hat{m}_h =$$

$$= m_g \bar{y}_{g0} - \frac{1}{G} \sum_{h=1}^G \bar{y}_{h0}$$

13.1.4 Within and between sums of squares, ANOVA F-test

110

Sums of squares

usual notation: $\bar{Y} = \frac{1}{n} \sum_{g=1}^G \sum_{j=1}^{n_g} Y_{g,j} = \frac{1}{n} \sum_{g=1}^G n_g \bar{Y}_g$

Within groups sum of squares (= residual sum of sq.)

$$SS_e = \|\hat{Y} - \hat{Y}\|^2 = \sum_{g=1}^G \sum_{j=1}^{n_g} |Y_{g,j} - \hat{Y}_{g,j}|^2 = \sum_{g=1}^G \sum_{j=1}^{n_g} |Y_{g,j} - \bar{Y}_g|^2$$

$$\nu_e = n - G$$

Between groups sum of squares (= regression sum of sq.)

$$SS_R = \|\hat{Y} - \bar{Y}\|^2 = \sum_{g=1}^G \sum_{j=1}^{n_g} |\hat{Y}_{g,j} - \bar{Y}|^2 = \sum_{g=1}^G n_g |\bar{Y}_g - \bar{Y}|^2$$

$$\nu_R = G - 1$$

One-way ANOVA F-test

111

consider $M: Y|Z \sim N(X\beta, \sigma^2 I_n)$

arbitrary parametrization of one-way ANOVA

$M_0: \sim N(\beta_0, \sigma^2 I_n)$

$M_0 \subset M$

$$SS_e^0 = SS_T = \sum_{g=1}^G \sum_{j=1}^{n_g} |Y_{g,j} - \bar{Y}|^2$$

$$SS_T = \underbrace{SS_e^0}_{SS_e} + SS_R$$

$$\Rightarrow SS_R = SS_e^0 - SS_e$$

\Rightarrow F-statistic to test $M_0 \subset M$:

$$F = \frac{\left(\frac{SS_R}{G-1}\right)}{\left(\frac{SS_e}{n-G}\right)}$$

MS_R : between variability

MSE : within variability

= classical ANOVA F-statistic

→ One-way ANOVA table

Effect (Term)	Degrees of Freedom	Effect sum of squares	Effect mean square	F-stat.	P-value
Factor	$G-1$	SS_R (between)	MS_R	F	p
Residual	$n-G$	SSE (within)	MSE		

