

5. ANOVA (Analysis of Variance)

[ANOVA 1] One-way ANOVA \Rightarrow Fixed effect regression model

$$(1) Y_{ij} = \mu_i + \epsilon_{ij} \quad i=1 \cdots K, \quad j=1, \dots, n_i$$

$$(2) Y_{ij} = \mu + \alpha_i + \epsilon_{ij} \quad \sum_{i=1}^K \alpha_i = 0$$

[ANOVA 2] (2.1) One-way ANOVA global F-test

Under $H_0: \mu_1 = \mu_2 = \dots = \mu_K \Leftrightarrow H_0: \alpha_1 = \alpha_2 = \dots = \alpha_K$

$$F_{\text{stat}} = \frac{(RSS_H - RSS_{\text{Full}}) / (K-1)}{RSS_{\text{Full}} / (n-K)} \sim F_{K-1, n-K}$$

$$RSS_{\text{Full}} = \sum_{i=1}^K \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2 \quad \hat{\mu}_i = \bar{Y}_{i.}$$

$$RSS_H = \sum_{i=1}^K \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2 \quad \hat{\mu}_H = \bar{Y}_{..}$$

(2.2) ANOVA table

Source	Sum of Squares	Degrees of freedom
Between groups	$SS_{\text{Between}} = \sum_{i=1}^K n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2$	$K-1$
Within groups	$SS_{\text{Within}} = \sum_{i=1}^K \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i.})^2$	$n-K$
Total	$SS_{\text{Total}} = \sum_{i=1}^K \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{..})^2$	$n-1$

\triangle When the global F-test is reject \Rightarrow look for contrasts

[ANOVA 3] One contrast

① Hypothesis test : under $H_0 : \mu_i = \mu_j$

$$\frac{\hat{\mu}_i - \hat{\mu}_j}{\text{S.E.}(\hat{\mu}_i - \hat{\mu}_j)} \sim t_{df} \quad \text{with} \quad \text{S.E.}(\hat{\mu}_i - \hat{\mu}_j) = \hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$$
$$df = n - k$$

② Confidence interval for $\mu_i - \mu_j$

$$\hat{\mu}_i - \hat{\mu}_j \pm t_{df}^{(\alpha/2)} \times \text{S.E.}(\hat{\mu}_i - \hat{\mu}_j)$$

Δ K groups $\Rightarrow \binom{K}{2}$ pairwise comparisons \Rightarrow multiple contrasts

[ANOVA 4] Multiple contrasts (Simultaneous inference)

Given multiple hypotheses H_{01}, \dots, H_{0m} for parameters $\theta_1, \dots, \theta_m$

(4.1) Familywise Error Rate (FWER) is probability of rejecting at least one of H_{01}, \dots, H_{0m} when they are all true

(4.2) Simultaneous confidence intervals at $100(1-\alpha)\%$ level

are intervals (L_i, U_i) $i=1, \dots, m$ with

$$P[L_i \leq \theta_i \leq U_i \text{ for all } i=1, \dots, m] > 1-\alpha$$

(4.3) Bonferroni's correction Let p_i be p-value of H_{0i} .

① Hypothesis test :

$$V1: \text{reject } H_{0i} \text{ if } p_i < \frac{\alpha}{m}$$

V2: define $P_{i,adj} = \min \{ m \times p_i, 1 \}$, reject H_0 if $P_{i,adj} < \alpha$

② Simultaneous confidence intervals: change quantile to at level $\frac{\alpha}{2m}$

E.g. $\hat{\mu}_i - \hat{\mu}_j \pm t_{df}^{(\frac{\alpha}{2m})} \times \text{s.e.}(\hat{\mu}_i - \hat{\mu}_j)$ with $m = \binom{K}{2}$, $df = n - k$

Differences and connections between ① one-way ANOVA F-test

② one contrast t-test ③ multiple contrast tests

- △ Bonferroni's correction is a general method.
- △ But it tends to be too conservative if $k > 10$.
- △ In specific problem, we may be able to improve.
- △ Suppose now we are interested in ALL pairwise comparisons $\mu_i = \mu_j$ for $i \neq j$

(4.4) Tukey - Kramer procedure for pairwise comparison

△ Studentized range distribution

- Let X_1, \dots, X_m be iid $N(\mu, \sigma^2)$
- Let $R = \max_{1 \leq i \leq m} X_i - \min_{1 \leq i \leq m} X_i$ be the range
- $\frac{R}{\hat{\sigma}}$ follows the studentized range distribution $q_{m, v}$
where m is # of normal r.v.

v is # of degrees of freedom used in estimating σ .

We apply the property to pairwise difference

△ Consider a balanced design $n_1 = \dots = n_k = n_B$

k groups have a common mean μ_0

$y_{ij} \sim N(\mu_0, \sigma^2)$ for $i=1 \dots k, j=1 \dots n_B$

$$|t_{ij}| = \frac{|\hat{\mu}_i - \hat{\mu}_j|}{\hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} = \frac{|\bar{y}_{i\cdot} - \bar{y}_{j\cdot}|}{\hat{\sigma} \sqrt{\frac{1}{n_B} + \frac{1}{n_B}}}$$

$$= \frac{1}{\sqrt{2}} \frac{1}{\hat{\sigma}} |\sqrt{n_B} \bar{y}_{i\cdot} - \sqrt{n_B} \bar{y}_{j\cdot}|$$

$$= \frac{1}{\sqrt{2}} \times \frac{1}{\hat{\sigma}} |\sqrt{n_B} (\bar{y}_{i\cdot} - \mu_0) - \sqrt{n_B} (\bar{y}_{j\cdot} - \mu_0)|$$

$$= \frac{1}{\sqrt{2}} \times \frac{1}{\hat{\sigma}} |x_i - x_j| \quad (*)$$

where we denote $x_i = \sqrt{n_B} (\bar{y}_{i\cdot} - \mu_0)$

- As $x_i \stackrel{iid}{\sim} N(0, \sigma^2)$, for $i=1 \dots k$

$q_{k, n-k}$ is the distribution of

$$\frac{\max_{1 \leq i \leq k} x_i - \min_{1 \leq i \leq k} x_i}{\hat{\sigma}}$$

$$= \max_{1 \leq i \neq j \leq k} (x_i - x_j) \times \frac{1}{\hat{\sigma}} \quad \text{By equivalence of definitions}$$

$$= \max_{1 \leq i < j \leq k} |x_i - x_j| \times \frac{1}{\hat{\sigma}} \Rightarrow \text{By } |x_1 - x_2| = \max\{x_1 - x_2, x_2 - x_1\}$$

$$= \sqrt{2} \times \max_{1 \leq i < j \leq k} |t_{ij}| \quad (\text{By } \textcircled{*} \text{ above})$$

$$\text{In summary } \sqrt{2} \times \max_{1 \leq i < j \leq k} |t_{ij}| \sim q_{k, n-k}$$

△ Let $q_{k, n-k}^{(\alpha)}$ denote the upper α -level quantile of the studentized range distribution

▷ [Hypothesis test with FWER control]

Reject $H_{0,i,j} : \mu_i = \mu_j$ if

$$|t_{ij}| = \frac{|\bar{y}_i - \bar{y}_j|}{\hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} > \frac{1}{\sqrt{2}} q_{k, n-k}^{(\alpha)}$$

- It controls FWER exactly at α for

balanced design $n_1 = \dots = n_k$

approximately at α for unbalanced design

Proof of FWER = α under balanced design

FWER

$$= P\left(\text{at least one of } |t_{ij}| > \frac{1}{\sqrt{2}} q_{k, n-k}^{(\alpha)} \mid \mu_i = \mu_j \text{ for all } i \neq j\right)$$

$$= P\left[\max_{1 \leq i < j \leq k} |t_{ij}| > \frac{1}{\sqrt{2}} q_{k, n-k}^{(\alpha)} \mid \mu_i = \mu_j \text{ for all } i \neq j\right]$$

$$= \alpha \quad \left(\begin{array}{l} \text{Because we have shown} \\ \sqrt{2} \times \max_{1 \leq i < j \leq k} |t_{ij}| \sim q_{k, n-k} \end{array} \right)$$

△ [Simultaneous confidence intervals]

Tukey's HSD intervals

(Honest Significance Difference)

$$\bar{y}_{i\cdot} - \bar{y}_{j\cdot} \pm \frac{1}{\sqrt{2}} q_{k, n-k}^{(\alpha)} \times \hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$$

14.5) False Discovery Rate (FDR)

△ Perform m simultaneous hypothesis tests

Classify the null hypotheses based on truth and results

	H_0 rejected	H_0 not rejected	Total
H_0 true	FD (V)	TN	m_0
H_0 false	TD	FN	m_1
Total	D (R) Rejected	N	m

T/F = true/false

P/N = Discovery / Nondiscovery
(Rejected)

△ All quantities except m , D and N are unobserved

False Discovery Proportion

$$\text{FDP} = \frac{\# \text{ False Discoveries}}{\# \text{ Discoveries}} = \frac{V}{\max\{R, 1\}}$$

$$\text{where } \max\{R, 1\} = \begin{cases} R & \text{if } R \geq 1 \\ 0 & \text{if } R = 0 \end{cases}$$

False Discovery Rate

$$\text{FDR} = E(\text{FDP}) = E\left(\frac{V}{\max\{R, 1\}}\right)$$

△ Benjamini-Hochberg (BH) procedure (JRSS-B 1995)

① Rank the p-values from smallest to largest with

$$P_{(0)} = 0 \quad P_{(1)} = \min_i P_i, \quad P_{(2)}, \quad P_{(3)} \cdots P_{(m)}$$

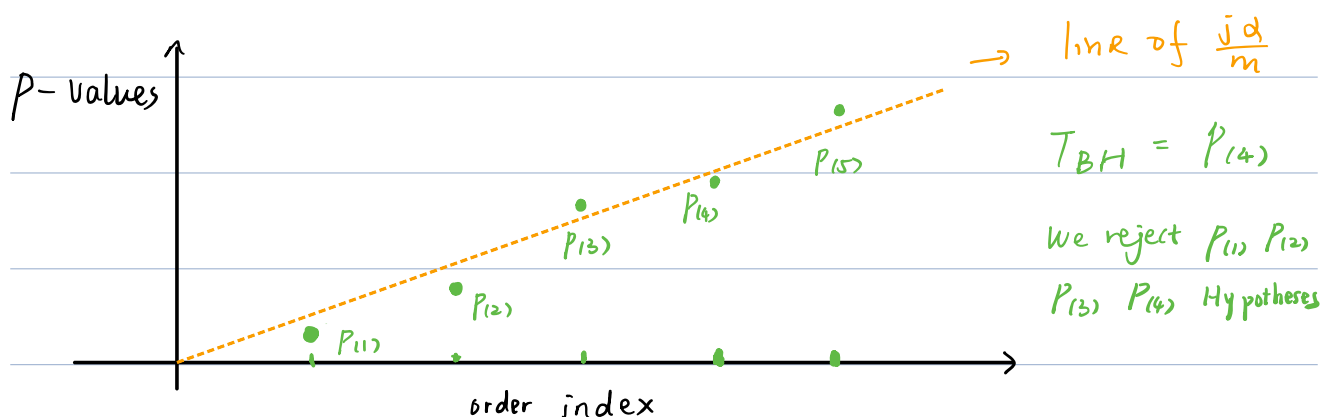
② Let the thresholds for $j = 0, \dots, m$ be

$$0, \quad \frac{\alpha}{m}, \quad \frac{2\alpha}{m}, \quad \frac{3\alpha}{m}, \quad \dots, \quad \alpha$$

③ BH threshold is defined for pre-specified $\alpha \in (0, 1)$

$$T_{\text{BH}} = \max_{1 \leq j \leq m} \left\{ P_{(j)} : P_{(j)} \leq \frac{j\alpha}{m} \right\}$$

Reject null hypothesis H_{0i} if $P_{(i)} \leq T_{\text{BH}}$



Equivalent procedure 1:

$\hat{j} = 4$ in above example

$$\text{Let } \hat{j} = \max_{1 \leq j \leq m} \left\{ j : P_{(j)} \leq \frac{j\alpha}{m} \right\}$$

Reject null hypotheses corresponding to $P_{(1)}, P_{(2)}, \dots, P_{(\hat{j})}$

Proof: Since $T_{BH} = P_{(\hat{j})}$ by definition,

$$P_{(l)} \leq T_{BH} \Leftrightarrow P_{(l)} \leq P_{(\hat{j})}$$

\Leftrightarrow By the order of $P_{(l)}$

l takes $\{1, 2, \dots, \hat{j}\}$

Equivalent procedure 2:

Define BH adjusted p-values

$$P_{(j), \text{adj}} = \min \left\{ \min_{l \geq j} \left\{ \frac{m P_{(l)}}{l} \right\}, 1 \right\}$$

Reject H_{0j} if $P_{(j), \text{adj}} \leq \alpha$

Proof: $P_{(j), \text{adj}} \leq \alpha \iff \min_{l \geq j} \frac{m P_{(l)}}{l} \leq \alpha$

Equivalent to saying $P_{(j)}$ is rejected if
one of index $l \geq j$ gives $P_{(l)} \leq \frac{l}{m} \alpha$

Looking at graphical example above

$P_{(3)}$ is rejected because $P_{(4)} \leq \frac{4}{m} \alpha$ (take $l=4$)

$P_{(4)}$ is rejected because $P_{(4)} \leq \frac{4}{m} \alpha$ (take $l=4$)

Theoretical guarantee

BH (1995) proved that for independent tests

the BH procedure guarantees that

$$\text{FDR} \leq \frac{m_0}{m} \alpha \leq \alpha$$

Intuition (only for students of interest, not required)

Suppose H_{0i} is true with probability $p = \frac{m_0}{m}$ independently

Reject at level t_{BH}

$$FDR = E\left(\frac{FD}{D}\right) \approx \frac{E(FD)/m}{E(D)/m}$$

$$= \frac{E\left\{\frac{1}{m} \sum_{i=1}^m \mathbb{1}(P_i \leq t_{BH}) \times \mathbb{1}(H_i \text{ holds})\right\}}{E\left\{\frac{1}{m} \sum_{i=1}^m \mathbb{1}(P_i \leq t_{BH})\right\}}$$

$$E\left\{\frac{1}{m} \sum_{i=1}^m \mathbb{1}(P_i \leq t_{BH})\right\}$$

$$\approx \frac{P(P_i \leq t_{BH} \mid H_i \text{ true}) P(H_i \text{ is true})}{\hat{G}(t_{BH})}$$

(Let $\hat{G}(t) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}(P_i \leq t)$ denote empirical CDF of P_i 's)

$$\approx \frac{t_{BH} \times \frac{m_0}{m}}{\hat{G}(t_{BH})}$$



$$\left[\begin{array}{l} \text{Use } P(P_i \leq t_{BH} \mid H_i \text{ true}) = t_{BH} \\ P(H_i \text{ is true}) = \frac{m_0}{m} \end{array} \right]$$

We want t_{BH} to be as large as possible while $\odot \leq \frac{m_0}{m} \alpha$

$$\hat{t}_{BH} = \sup \left\{ t : \frac{t \times \frac{m_0}{m}}{\hat{G}(t)} \leq \frac{m_0}{m} \alpha \right\}$$

$$\approx \max_{1 \leq j \leq m} \left\{ t_{(j)} : t_{(j)} \leq \hat{G}(t_{(j)}) \times \alpha \right\}$$

(Since $t_{(j)}$ j -th smallest p -value, $\hat{G}(t_{(j)}) = j/m$)

$$\approx \max_{1 \leq j \leq m} \left\{ t_{(j)} : t_{(j)} \leq \frac{j}{m} \times \alpha \right\}$$