

ANOVA

Summary

Can be applied to **NESTED** model selection

Analysis Of Variance

one-way

[E.g.] IQ of different schools (单因素)

Testing: $H_0: \mu_1 = \mu_2 = \dots = \mu_k$ $H_A: H_0$ is false

School 1	$IQ_{11} \dots$	IQ_{1n_1}	n_1	\overline{IQ}_1
School 2	$IQ_{21} \dots$	IQ_{2n_2}	n_2	\overline{IQ}_2
\vdots	\vdots	\vdots	\vdots	\vdots
School K	IQ_{K1}	IQ_{Kn_k}	n_k	\overline{IQ}_k

\overline{IQ}

two-way

IQ of different schools & sex (双因素)

(Similar to 单因素 scenario)

school 1
 \vdots
 school K

Male
 Female

(F-test)

① Understanding (One-way) from LR perspective

Condition

a) suppose G-M condition

① $E[\epsilon_i] = 0$
 ② $cov(\epsilon_i, \epsilon_j) = 0$

b) $\epsilon \sim \text{Gaussian}$

Problem Setting:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_{p-1} X_{p-1,i} + \epsilon_i$$

$$X_{1i}, \dots, X_{pi} \sim 0/1$$

→ avoid col-linearity issue

dummy variable

⇔ group 1
⋮
group p

Since here we only have $p-1$ variables!

★ \hat{y}_i is derived from full model

$$SSR = \sum (\hat{y}_i - \bar{y})^2$$

$$SSE = \sum (y_i - \hat{y}_i)^2$$

LR Framework

$$\frac{\frac{SSR}{p-1}}{\frac{SSE}{n-p}} \sim F$$

→ Q: Here, what is \hat{y}_i ?

Logic: ① we achieve SSR/SSE under full model

② $\frac{\frac{SSR}{p}}{\frac{SSE}{n-p+1}} \sim F$ under null assumption

→ Answer: ANOVA Test ⇔ $\begin{cases} H_0: \beta_1 = \dots = \beta_{p-1} = 0 \\ H_A: \text{o/w} \end{cases}$

★ $\frac{\frac{SSR}{p-1}}{\frac{SSE}{n-p}} \sim F$

→ ① Notice that, under H_0 , model turns to

★ $\Rightarrow \hat{\beta}_0 = \bar{y}$

$y_i = \beta_0 + \varepsilon_i$

→ ② Try to figure out \hat{y}_i !

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{p-1,i} + \varepsilon_i$$

$$\rightarrow \hat{\beta} = \min_{\beta} \| \hat{y} - y \|_2^2$$

$$= \min_{\beta} \sum_{y \in G_1} (y - \beta_0 - \beta_1)^2 + \dots + \sum_{y \in G_{p-1}} (y - \beta_0 - \beta_{p-1})^2 + \sum_{y \in G_p} (y - \beta_0)^2$$

$$\Rightarrow \begin{cases} \hat{\beta}_0 = \bar{y}_p \\ \hat{\beta}_0 + \hat{\beta}_1 = \bar{y}_1 \\ \vdots \\ \hat{\beta}_0 + \hat{\beta}_{p-1} = \bar{y}_{p-1} \end{cases}$$

⇒

$$\begin{cases} \hat{y}_i = \bar{y}_{g(i)} \\ g(i) = \{ \text{Group: } y_i \text{ belongs to} \} \end{cases}$$

\Rightarrow Therefore, we can compute $\begin{cases} SSR \\ SSE \end{cases}$

$$\textcircled{1} SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$= \sum_{k=1}^p (\# \text{ of data in Group } k) \cdot (\bar{y}_k - \bar{y})^2$$

$$\textcircled{2} SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{k=1}^p \sum_{y_i \in \text{Group } k} (y_i - \bar{y}_k)^2$$

$$\Rightarrow \frac{\frac{SSR}{p-1}}{\frac{SSE}{n-p}} \sim F(p-1, n-p)$$

(组间) (组内)

Another Perspective: Analyze Directly (Common Manner)

\rightarrow From 学弱猫

$\textcircled{1}$ Model Setting

Notation

$\begin{cases} \mu \rightarrow \text{mean of all group} \\ a_i \rightarrow i\text{-th group power} \end{cases}$

$$\begin{cases} y_{ij} = \mu + a_i + \varepsilon_{ij} & i=1, 2, \dots, r ; j=1, 2, \dots, m \\ \sum_{i=1}^r a_i = 0 \\ \mu_i = \mu + a_i & (y_{ij} = \mu_i + \varepsilon_{ij}) \\ \varepsilon_{ij} \sim N(0, \sigma^2) & \& \text{ independent with each other} \end{cases}$$

ANOVA \Leftrightarrow Hypothesis Testing on $\begin{cases} H_0: a_1 = a_2 = \dots = a_r = 0 \\ H_a: \text{otherwise} \end{cases}$

② Analysis:

$$\begin{cases} \bar{y}_i = \mu + a_i + \bar{\varepsilon}_i \\ \bar{y} = \mu + \bar{\varepsilon} \end{cases}$$

$$SSR = \sum_{i=1}^r m \cdot (\bar{y}_i - \bar{y})^2$$

$$= \sum_{i=1}^r m (\bar{\varepsilon}_i + a_i - \bar{\varepsilon})^2$$

$$\underline{\underline{H_0}} \quad \sum_{i=1}^r m (\bar{\varepsilon}_i - \bar{\varepsilon})^2 \sim \sigma^2 \chi^2(r-1)$$

$$SSE = \sum_{i=1}^r \sum_{j=1}^m (y_{ij} - \bar{y}_i)^2$$

$$= \sigma^2 (\chi_1^2(m-1) + \dots + \chi_r^2(m-1))$$

$$= \sigma^2 \chi^2(r \cdot m - r)$$

Notice that $\begin{cases} (y_{ij} - \bar{y}_i)^2 \perp \bar{y}_i \\ (y_{ij} - \bar{y}_i)^2 \perp \bar{y}_\ell \quad \ell \neq i \end{cases} \Rightarrow (y_{ij} - \bar{y}_i)^2 \perp (\bar{y}_\ell - \bar{y})^2 \quad \forall \ell$
 $\Rightarrow \underline{\underline{SSR \perp SSE}}$

$$\Rightarrow \frac{\frac{SSR}{r-1}}{\frac{SSE}{r(m-1)}} \sim F(r-1, r(m-1))$$

ANOVA Table

	value	df	Normalize	F	P
SSR	$\sum_{i=1}^r m (\bar{y}_i - \bar{y})^2$	$r-1$	$SSR/r-1$	$\frac{SSR}{r-1}$	$P(F > F_0)$
SSE	$\sum_{i=1}^r \sum_{j=1}^m (y_{ij} - \bar{y}_i)^2$	$r(m-1)$	$SSE/r(m-1)$	$\frac{SSE}{r(m-1)}$	

To compute efficiently, we always compute

$$SST = \sum_{i=1}^r \sum_{j=1}^m (y_{ij} - \bar{y})^2$$

$$= \sum_{i=1}^r \sum_{j=1}^m y_{ij}^2 - \frac{T^2}{n}$$

$$SSR = \sum_{i=1}^r m (\bar{y}_i - \bar{y})^2$$

$$= \sum_{i=1}^r m \left(\frac{T_i}{m} - \frac{T}{n} \right)^2$$

$$= \sum_{i=1}^r \frac{T_i^2}{m} - 2 \cdot \frac{T_i \cdot T}{n} + m \cdot \frac{T^2}{n^2}$$

$$= \frac{1}{m} \sum_{i=1}^r T_i^2 - \frac{T^2}{n}$$

$$T = \sum_{i,j} y_{ij} \quad n = m \cdot r$$

$$T^2 = \sum_{i,j} y_{ij}^2$$

$$T_i = \sum_j y_{ij}$$

ANOVA → actually one special realization of F-test
in conventional LR model

Terminology: ① $[X_1, \dots, X_p]$ is Continuous variables

test $p \in S, \beta_p = 0 \Rightarrow$ **F-test**

☆

$$\frac{\frac{SSR}{p}}{\frac{SSE}{n-p-1}} \sim F$$

is one specific form

Constraint LR Model

↓

$$\frac{\frac{SSE(H_0) - SSE}{df_1}}{\frac{SSE}{n-p-1}} \sim F$$

② $\text{Null} \subseteq A \subseteq A + A:B \dots$ \rightsquigarrow Nested Factors

$A, B \rightsquigarrow$ Factors (discrete variables)

test: $\beta_{A1} = \beta_{A2} = \dots = \beta_{AP_A} = 0$

($\beta_{B11} = \dots = \beta_{B1P_B} = \dots = \beta_{BP_AP_B} = 0$)

\rightarrow principally, apply F-test in LR Framework!

ANOVA

Remark: 1.) ① and ② are actually equivalent, while ② have the "nested" constraints!

2) From ②, we care more about the quantitative conclusion. That is, whether these factors really contribute to our predicted target (hypothesis testing result)

in ANOVA, this can be done without

actually going through parameters estimation

$$\text{RSS} \leftrightarrow \text{SSE} \\ (y_{ij} - \hat{y}_{ij})^2$$

One-way [M_1 vs M_2]

[E.g.]

$$\frac{\text{Null}}{M_1} \subseteq \frac{A}{M_2} \subseteq \frac{A + A:B}{M_3}$$

$$\text{RSS}(\text{Null}) - \text{RSS}(M_1) \\ = (\bar{y}_i - \bar{y})^2$$

Test $A:B$ significance

$$\text{SSR} = \text{RSS}(M_2) - \text{RSS}(M_3)$$

$$\text{SSE} = \text{RSS}(M_3)$$

