

第十八讲. 一般线性假设及单因素方差分析

1

回顾: 部分回归系数的显著性检验

全模型: $Y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \varepsilon_{n \times 1} = X_1 \beta_{(1)} + X_2 \beta_{(2)} + \varepsilon, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 I_n)$

其中 $X = (X_1, X_2)$, X_1 第一列为 $\mathbf{1}$, $\beta = \begin{pmatrix} \beta_{(1)} \\ \beta_{(2)} \end{pmatrix}$, $\beta_{(2)}$ 长度为 q .

原假设 $H_0: \beta_{(2)} = \mathbf{0}_{q \times 1}$,

子模型: $Y = X_1 \beta_{(1)} + \varepsilon, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 I_n)$ (原假设下的模型)

$$F = \frac{\|\hat{Y} - \hat{Y}_0\|^2 / q}{\|Y - \hat{Y}\|^2 / (n-p)} = \frac{\|X_2^\perp \hat{\beta}_{(2)}\|^2 / q}{\hat{\sigma}^2}$$

其中 $\hat{Y} = P_X Y$, $\hat{Y}_0 = P_{X_1} Y$, 且 $\hat{Y} - \hat{Y}_0 = X_2^\perp \hat{\beta}_{(2)}$ 其中 $X_2^\perp = X_2 - P_{X_1} X_2$

2

定理: 原假设下 $F \sim F_{q, n-p}$

证明1: 利用 $F = \frac{\|X_2^\perp \hat{\beta}_{(2)}\|^2}{q \hat{\sigma}^2}$.

由 $\hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1})$, 知 $\hat{\beta}_{(2)} \sim N(\beta_{(2)}, \sigma^2 (X_2^\perp{}' X_2^\perp)^{-1})$

H_0 成立时, $\beta_{(2)} = 0$, 所以 $A = \|X_2^\perp \hat{\beta}_{(2)}\|^2 = \hat{\beta}_{(2)}' X_2^\perp{}' X_2^\perp \hat{\beta}_{(2)} \sim \sigma^2 \chi_q^2$

另外, $B = (n-p) \hat{\sigma}^2 \sim \sigma^2 \chi_{n-p}^2$, 且与 $\hat{\beta}$ 独立,

所以 $\frac{A/q}{B/(n-p)} = \frac{\|X_2^\perp \hat{\beta}_{(2)}\|^2}{q \hat{\sigma}^2} = F \sim F_{q, n-p}$

3

证明2: 利用 $F = \frac{n-p}{q} \times \frac{\|\hat{Y} - \hat{Y}_0\|^2}{\|Y - \hat{Y}\|^2}$.

(1) $Y - \hat{Y} = (I_n - P_X)Y = (I_n - P_X)\varepsilon$

$\frac{(n-p)\hat{\sigma}^2}{\sigma^2} = \frac{\|Y - \hat{Y}\|^2}{\sigma^2} \sim \chi_{n-p}^2$ (不论 H_0 成立与否).

(2) 原假设下 $Y = X_1 \beta_{(1)} + \varepsilon, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 I_n)$

$\Rightarrow \hat{Y} - \hat{Y}_0 = P_{X_2^\perp} Y = P_{X_2^\perp} (X_1 \beta_{(1)} + \varepsilon) = P_{X_2^\perp} \varepsilon$

$\Rightarrow \|\hat{Y} - \hat{Y}_0\|^2 = \varepsilon' P_{X_2^\perp} \varepsilon \sim \sigma^2 \chi_q^2$

(3) 因为 $(I_n - P_X)P_{X_2^\perp} = 0$, 所以 $\|\hat{Y} - \hat{Y}_0\|^2$ 与 $\|Y - \hat{Y}\|^2$ 独立,

$\Rightarrow F = \frac{\|\hat{Y} - \hat{Y}_0\|^2 / q}{\|Y - \hat{Y}\|^2 / (n-p)} \sim F_{q, n-p}$

4

ANOVA: F可以表达为比较模型的拟合值、残差、拟合优度

$$F = \frac{\|\hat{Y} - \hat{Y}_0\|^2 / q}{\|Y - \hat{Y}\|^2 / (n-p)}$$

$$F = \frac{n-p}{q} \times \frac{RSS_0 - RSS}{RSS}$$

$$F = \frac{n-p}{q} \times \frac{R^2 - R_0^2}{1 - R^2}$$

5

回归方程的显著性检验

$$Y = \mathbf{1}\beta_0 + Z\gamma + \varepsilon, \beta_0 \text{ 是截距}$$

$$H_0: \gamma = (\beta_1, \dots, \beta_{p-1})' = 0$$

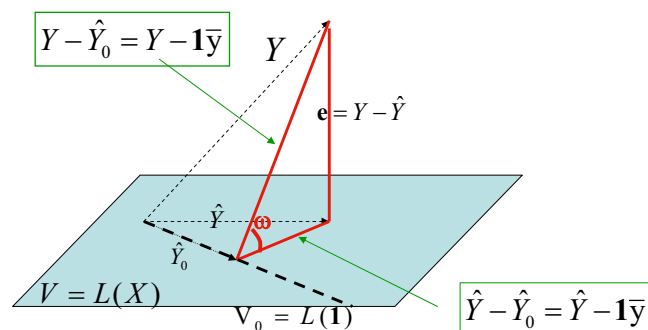
记 $Z^\perp = Z - P_1 Z = Z - \mathbf{1}\bar{x}'$ 为自变量的中心化矩阵。

回归方程显著性的 F 检验:

$$F = \frac{n-p}{p-1} \times \frac{\|\hat{Y} - \hat{Y}_0\|^2}{\|Y - \hat{Y}\|^2} = \frac{\|Z^\perp \hat{\gamma}\|^2}{(p-1)\hat{\sigma}^2} \sim_{H_0} F_{p-1, n-p}$$

其中 $\hat{Y}_0 = \mathbf{1}\bar{y}$, 所以 $F = \frac{n-p}{p-1} \times \frac{R^2}{1-R^2}$

6



$$F = \frac{n-p}{p-1} \times \frac{\|\hat{Y} - \hat{Y}_0\|^2}{\|Y - \hat{Y}\|^2} = \frac{n-p}{p-1} \times \frac{R^2}{1-R^2} = \frac{n-p}{p-1} \times (\cot \omega)^2$$

$$\cos^2 \omega = \frac{\|\hat{Y} - \mathbf{1}\bar{y}\|^2}{\|Y - \mathbf{1}\bar{y}\|^2} = R^2, \omega = \text{中心化后 } Y \text{ 与 } X \text{ 平面的最大夹角}$$

一般线性假设

一般线性假设通常表述为

$$H_0: A\beta = \mathbf{0}_{q \times 1}, q \leq p, A \text{ 为 } q \times p \text{ 已知矩阵 (行满秩)}.$$

记均值向量 $\mu = X\beta$, 全模型可表示为: $Y = \mu + \varepsilon, \mu \in V = L(X)$

原假设下 β 处于 q 个约束之下, 则 $\mu \in$ 某个 $V_0 \subset V, V_0$ 由 A 决定 (比如若 $A = (0, I_q)$, 则 $A\beta = \beta_{(2)} = 0, \mu = X\beta = X_1\beta_{(1)} \in L(X_1) \subset V$).

$$\text{记 } \hat{Y} = P_V Y, \hat{Y}_0 = P_{V_0} Y$$

一般线性假设 ($H_0: \mu \in V_0$) 的 F-检验为:

$$F = \frac{\|\hat{Y} - \hat{Y}_0\|^2 / q}{\|Y - \hat{Y}\|^2 / (n-p)} = \frac{(RSS_0 - RSS)/q}{RSS/(n-p)} \sim_{H_0} F_{q, n-p}$$

原假设下的 A, X 决定的投影空间 V_0 是什么?

8

$$\text{例1. } H_0: A\beta = 0, \text{ 其中 } A = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 \\ & & & \cdots & & & \\ 0 & 0 & 0 & 0 & \cdots & 1 & -1 \end{pmatrix}$$

$$\Leftrightarrow H_0: \beta_1 = \dots = \beta_{p-1},$$

H_0 成立时

$$\begin{aligned} Y &= \mathbf{1}\beta_0 + \mathbf{x}_1\beta_1 + \dots + \mathbf{x}_{p-3}\beta_{p-3} + \mathbf{x}_{p-2}\beta_{p-2} + \mathbf{x}_{p-1}\beta_{p-1} + \varepsilon \\ &= \mathbf{1}\beta_0 + (\mathbf{x}_1 + \dots + \mathbf{x}_{p-1})\gamma + \varepsilon \end{aligned}$$

记 $X_0 = (\mathbf{1}, \mathbf{x}_1 + \dots + \mathbf{x}_{p-1})$, 所以 $V_0 = L(X_0)$

$$\text{例2. } A = \begin{pmatrix} 0 & 1 & 1 & -2 & 0 & \dots & 0 \end{pmatrix}, H_0: A\beta = 0, \Leftrightarrow H_0: \beta_1 + \beta_{p-3} = 2\beta_3,$$

H_0 成立时

$$\begin{aligned} Y &= \mathbf{1}\beta_0 + \mathbf{x}_1\beta_1 + \mathbf{x}_2\beta_2 + \mathbf{x}_3\beta_3 + \mathbf{x}_4\beta_4 + \dots + \varepsilon \\ &= \mathbf{1}\beta_0 + (\mathbf{x}_1 + \mathbf{x}_3/2)\beta_1 + (\mathbf{x}_2 + \mathbf{x}_3/2)\beta_2 + \mathbf{x}_4\beta_4 + \dots + \varepsilon, \end{aligned}$$

记 $X_0 = (\mathbf{1}, \mathbf{x}_1 + \mathbf{x}_3/2, \mathbf{x}_2 + \mathbf{x}_3/2, \mathbf{x}_4, \dots, \mathbf{x}_{p-1})$, 所以 $V_0 = L(X_0)$

9

对于一般线性假设:

$$H_0: A\beta = \mathbf{0}_{q \times 1}, \quad q \leq p, A \text{ 为 } q \times p \text{ 已知矩阵(行满秩).}$$

由 $\mu = X\beta \Rightarrow \beta = (X'X)^{-1}X'\mu$, 所以

$$A\beta = A(X'X)^{-1}X'\mu = 0 \Leftrightarrow \mu \perp L(X(X'X)^{-1}A'), \quad \mu \in L(X)$$

$$V_0 = L(X(X'X)^{-1}A')^\perp \cap L(X) \Rightarrow P_{V_0} = P_X - P_{X(X'X)^{-1}A'},$$

$$\hat{Y}_0 = \hat{Y} - X(X'X)^{-1}A'[A(X'X)^{-1}A']^{-1}A(X'X)^{-1}X'Y$$

$$\Rightarrow \|\hat{Y} - \hat{Y}_0\|^2 = \hat{\beta}'A'[A(X'X)^{-1}A']^{-1}A\hat{\beta} \Rightarrow$$

$$F = \frac{n-p}{q} \frac{\|\hat{Y} - \hat{Y}_0\|^2}{\|Y - \hat{Y}\|^2} = \hat{\beta}'A'[A(X'X)^{-1}A']^{-1}A\hat{\beta}/q\hat{\sigma}^2$$

10

注1: 事实上, $H_0: A\beta = 0$ 的检验可直接构造得到:

$$H_0 \text{ 下, } A\hat{\beta} \sim N_q(0, \sigma^2 A(X'X)^{-1}A') \Rightarrow \hat{\beta}'A'(A(X'X)^{-1}A')^{-1}A\hat{\beta}/\sigma^2 \sim \chi_q^2$$

度量了 H_0 成立的证据, plug-in $\hat{\sigma}^2$ 并除以 q , 即得到

$$F = \hat{\beta}'A'(A(X'X)^{-1}A')^{-1}A\hat{\beta}/q\hat{\sigma}^2 \stackrel{H_0}{\sim} F_{q, n-p}$$

注2: 如果原假设具有如下形式

$$H_0: A\beta = \mathbf{c}_{q \times 1} \quad (\mathbf{c} \text{ 已知}), A \text{ 为 } q \times p \text{ 已知矩阵(行满秩)}$$

$$F = \frac{(A\hat{\beta} - \mathbf{c})'(A(X'X)^{-1}A')^{-1}(A\hat{\beta} - \mathbf{c})}{q\hat{\sigma}^2} \sim_{H_0} F_{q, n-p}$$

比如, R输出中给出了 $\hat{\beta}$ 及其标准差 $se(\hat{\beta})$, 以及 $H_0: \beta = 0$ 的

t-检验。如果检验 $H_0: \beta = \beta_0$ (已知), 只需计算 $t = (\hat{\beta} - \beta_0)/se(\hat{\beta})$

11

单因素方差分析(one-way anova)

随机化控制试验: 某因子变量有 K 个水平(处理, treatment), 对研究对象随机分配处理(随机分组), 观察响应 y , 考察因子变量有无效应。数据和模型假设:

第 k 组 $y_{k1}, \dots, y_{kn_k} \sim N(\mu_k, \sigma^2), k = 1, \dots, K$, 各组独立,

$$H_0: \mu_1 = \mu_2 = \dots = \mu_K$$

12

$$\begin{array}{c}
 \text{第一组} \\
 \vdots \\
 y_{1n_1} \\
 \hline
 \text{第二组} \\
 y_{21} \\
 \vdots \\
 y_{2n_2} \\
 \vdots \\
 \hline
 \text{第 } K \text{ 组} \\
 y_{K1} \\
 \vdots \\
 y_{Kn_K}
 \end{array}
 =
 \begin{array}{c}
 \mathbf{x}_1 \quad \mathbf{x}_2 \dots \mathbf{x}_K \\
 \begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 1 \end{pmatrix}
 \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_K \end{pmatrix}
 + \boldsymbol{\varepsilon}
 \end{array}$$

13

全模型 $Y = \mathbf{x}_1\mu_1 + \mathbf{x}_2\mu_2 + \dots + \mathbf{x}_K\mu_K + \varepsilon$

$\hat{Y} = (\bar{y}_{1\cdot}, \dots, \bar{y}_{1\cdot}, \bar{y}_{2\cdot}, \dots, \bar{y}_{2\cdot}, \dots, \bar{y}_{K\cdot}, \dots, \bar{y}_{K\cdot})$,

其中 $\bar{y}_{k\cdot} = (y_{k1} + \dots + y_{kn_k}) / n_k$

原假设下 $\mu_1 = \mu_2 = \dots = \mu_K$,

$Y = (\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_K)\mu_1 + \varepsilon = \mathbf{1}\mu_1 + \varepsilon$

$\hat{Y}_0 = (\bar{y}_{\cdot\cdot}, \dots, \bar{y}_{\cdot\cdot})$, $\bar{y}_{\cdot\cdot} = \sum \sum y_{ij} / n$, $n = n_1 + \dots + n_K$

$$\begin{aligned}
 F &= \frac{\|\hat{Y} - \hat{Y}_0\|^2 / (K-1)}{\|Y - \hat{Y}\|^2 / (n-K)} = \frac{(n_1(\bar{y}_{1\cdot} - \bar{y}_{\cdot\cdot})^2 + \dots + n_K(\bar{y}_{K\cdot} - \bar{y}_{\cdot\cdot})^2) / (K-1)}{\sum_{k=1}^K \sum_{j=1}^{n_k} (y_{kj} - \bar{y}_{k\cdot})^2 / (n-K)} \\
 &= \frac{SS_{\text{组间}} / (K-1)}{SS_{\text{组内}} / (n-K)} \sim_{H_0} F_{K-1, n-K}
 \end{aligned}$$

14

事实上，通常的做法是直接进行方差（平方和）分解：

$$\begin{aligned}
 SS_{\text{总}} &= \sum_{k=1}^K \sum_{j=1}^{n_k} (y_{kj} - \bar{y}_{\cdot\cdot})^2 = \sum_{k=1}^K \sum_{j=1}^{n_k} (y_{kj} - \bar{y}_{k\cdot} + \bar{y}_{k\cdot} - \bar{y}_{\cdot\cdot})^2 \\
 &= \sum_{k=1}^K \sum_{j=1}^{n_k} (\bar{y}_{k\cdot} - \bar{y}_{\cdot\cdot})^2 + \sum_{k=1}^K \sum_{j=1}^{n_k} (y_{kj} - \bar{y}_{k\cdot})^2 \triangleq SS_{\text{组间}} + SS_{\text{组内}} \\
 F &= \frac{SS_{\text{组间}} / (K-1)}{SS_{\text{组内}} / (n-K)} \sim F_{K-1, n-K} \quad (\text{原假设下})
 \end{aligned}$$

称为单因素方差分析(one-way anova)

方差分析表：

来源	因子/分组	残差
平方和	$SS_{\text{组间}}$	$SS_{\text{组内}}$
自由度	$K-1$	$n-k$

15

AOV函数(Analysis of Variance)

> aov(y ~ group)

group 代表了每个观察所属组别，是因子变量。

> aov(TS~D, data=sleep1)

Call:

aov(formula = TS ~ D, data = sleep1)

Terms:

	D	Residuals
Sum of Squares	457.2556	752.4122
Deg. of Freedom	4	53

← 方差分析表

Residual standard error: 3.767818

Estimated effects may be unbalanced

4 observations deleted due to missingness

F= (457.2556/4)/(752.4122/ 53)=8.052

16

非参数检验方法（不假设正态分布）

两样本: Wilcoxon秩和检验（或Mann-Whitney- Wilcoxon检验）

处理组: y_1, \dots, y_{n_1} iid $\sim f(x - \mu_1)$, f 密度, 但未知, μ_1 为均值

对照组: $y_{n_1+1}, \dots, y_{n_1+n_2}$ iid $\sim f(x - \mu_2)$, μ_2 为均值

$H_0: \mu_1 = \mu_2$

```
> x=c(12,22,8,5,30)
> rank(x)
[1] 3 4 2 1 5
```

y_i 在所有 $n = n_1 + n_2$ 个观察值中的排名/秩记为 R_i

第一组的平均秩 $\bar{R} = (R_1 + \dots + R_{n_1}) / n_1$

第二组的平均秩 $\bar{R} = (R_{n_1+1} + \dots + R_{n_1+n_2}) / n_2$

17

原假设 $H_0: \mu_1 = \mu_2$ 下, 两组的平均秩应该相差不大,

Wilcoxon检验统计量为:

$$W = \bar{R}_2 - \bar{R}_1$$

如何评价 W 的显著性? 只要能否算出其 p 值.

原假设 $H_0: \mu_1 = \mu_2$ 下, 两组无差异, 那么我们随机置换 $y_1, \dots, y_{n_1}, \dots, y_{n_1+n_2}$, 即随机分组, 置换后的数据重新计算平均秩的差值, 记为 W^* , 如果 H_0 成立, 那么 $W^* \approx W$

反复置换 N 次, 所得平均秩之差为 W_1^*, \dots, W_N^*

W 的 p 值: $p = \frac{\{ |W_i^*| > |W| \text{ 的个数} \}}{N}$

```
> wilcox.test ( y ~ group)
```

18

K样本: Kruskal-Wallis检验

组1: y_1, \dots, y_{n_1} iid $\sim f(x - \mu_1)$, f 密度, 但未知, μ_1 为均值

组2: $y_{n_1+1}, \dots, y_{n_1+n_2}$ iid $\sim f(x - \mu_2)$,

...

组 K : y_{n-n_K+1}, \dots, y_n iid $\sim f(x - \mu_K)$,

$H_0: \mu_1 = \mu_2 = \dots = \mu_K$

y_i 在所有 $n = n_1 + n_2$ 个观察值中的排名/秩记为 R_i

记第 k 组的平均秩 \bar{R}_k , 原假设下 $\bar{R}_1, \dots, \bar{R}_K$ 应该差别不大.

```
> kruskal.test ( y ~ group)
```

19