

第七章 方差分析模型

方差分析模型是应用非常广泛的一类线性模型. 这种模型多有一定的试验设计背景, 因而也被称为试验设计模型. 方差分析由英国著名的统计学家Fisher于20世纪20年代提出, 有时也称为 F 检验. 从形式上看, 方差分析是用来比较多个总体均值的假设检验方法, 其本质是研究一个(或多个)定性变量(即分类变量)与定量变量之间的关系.

方差分析的应用范围很广, 可用于社会科学、生物工程、医学研究等各个领域的试验数据分析. 本章仅仅介绍最基本的单因素方差分析和两因素方差分析.

7.1 单因素方差分析

方差分析起源于农业田间试验. 假定某个农业试验基地引进了 a 个小麦品种, 在进行大面积种植之前, 先进行小范围的试验种植, 以便从中挑选出最适合本地区的品种.

将一大块田划分成面积相等的几个小块, 其中 n_1 块种植第1种小麦, n_2 块种植第2种小麦, \dots , 等等. 试验的目的是比较小麦的品种, 因此我们感兴趣的只是小麦品种这一个因素. 其它所有因素, 如施肥量、浇水等对这 n 块田都控制在相同状态下. 在这个例子里, 我们感兴趣的因素只有一个, 即小麦品种. 每个具体的品种, 都称为小麦品种这个因素的一个水平. 现有 a 个不同的品种, 因此小麦品种这一因素一共有 a 个水平. 这是单因素 a 个水平的问题.

记 y_{ij} 为种植第 i 个品种的小麦在第 j 块田的产量, $i = 1, \dots, a; j = 1, \dots, n_i$. 对固定的 i , y_{i1}, \dots, y_{in_i} 分别是第 i 个品种的小麦在 n_i 块田的产量. 因为除了一些随机误差外, 这 n_i 块田的一切生产条件完全一样. 因此可把它们看做来自正态总体的样本. 即, 可假设 $\{y_{ij}, i = 1, \dots, a; j = 1, \dots, n_i\}$ 相互独立且

$$y_{ij} \sim N(\mu_i, \sigma^2), \quad i = 1, \dots, a; \quad j = 1, \dots, n_i. \quad (7.1.1)$$

表7.1.1描述了此问题的总体和样本.

水平	总体	样本
1	$N(\mu_1, \sigma^2)$	$y_{11}, y_{12}, \dots, y_{1n_1}$
2	$N(\mu_2, \sigma^2)$	$y_{21}, y_{22}, \dots, y_{2n_2}$
\vdots	\vdots	\vdots
a	$N(\mu_a, \sigma^2)$	$y_{a1}, y_{a2}, \dots, y_{an_a}$

表7.1.1 单因素方差分析问题

称所考虑的因素为因素 A , 假定它有 a 个水平, 我们的目的是比较这 a 个水平的差异. 将(7.1.1)改写为

$$\begin{cases} y_{ij} = \mu_i + e_{ij}, \\ e_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2), \end{cases} \quad i = 1, \dots, a; \quad j = 1, \dots, n_i, \quad (7.1.2)$$

其中 e_{ij} 是试验误差. 容易看出, 比较因素 A 的 a 个水平的差异可归结为比较这 a 个总体的均值 μ_1, \dots, μ_a 之间的差异.

记

$$\mu = \frac{1}{n} \sum_{i=1}^a n_i \mu_i, \quad n = \sum_{i=1}^a n_i, \quad \alpha_i = \mu_i - \mu,$$

这里, μ 为整个样本的均值的总平均, α_i 表示第 i 个水平下的均值与总平均的差异, 它反映了第 i 个水平对指标 y 的效应. 因此有

$$\sum_{i=1}^a n_i \alpha_i = \sum_{i=1}^a n_i (\mu_i - \mu) = n\mu - n\mu = 0.$$

因为 $\mu_i = \mu + \alpha_i$, 于是(7.1.2)又可以写为

$$\begin{cases} y_{ij} = \mu + \alpha_i + e_{ij}, \\ e_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2), \quad i = 1, \dots, a; \quad j = 1, \dots, n_i. \\ \sum_{i=1}^a n_i \alpha_i = 0, \end{cases} \quad (7.1.3)$$

这就是单因素方差分析模型. 写成矩阵形式即为

$$\begin{cases} \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n), \\ \mathbf{h}'\boldsymbol{\beta} = 0, \end{cases} \quad (7.1.4)$$

其中

$$\begin{aligned} \mathbf{Y} &= (y_{11}, \dots, y_{1n_1}, y_{21}, \dots, y_{2n_2}, \dots, y_{a1}, \dots, y_{an_a})', \\ \boldsymbol{\beta} &= (\mu, \alpha_1, \dots, \alpha_a)', \\ \mathbf{e} &= (e_{11}, \dots, e_{1n_1}, e_{21}, \dots, e_{2n_2}, \dots, e_{a1}, \dots, e_{an_a})', \\ \mathbf{h} &= (0, n_1, n_2, \dots, n_a)', \end{aligned}$$

以及

$$\mathbf{X} = \begin{matrix} \text{第1行} \\ \vdots \\ \text{第}n_1\text{行} \\ \text{第}n_1+1\text{行} \\ \vdots \\ \text{第}n_1+n_2\text{行} \\ \vdots \\ \text{第}n_1+\dots+n_{a-1}+1\text{行} \\ \vdots \\ \text{第}n_1+\dots+n_{a-1}+n_a\text{行} \end{matrix} \begin{pmatrix} 1 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \cdots & 1 \end{pmatrix}.$$

可见, 单因素方差分析模型是一个带约束条件($\mathbf{h}'\boldsymbol{\beta} = 0$)的线性模型.

可以看出, 检验

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_a$$

等价于检验

$$H_0 : \alpha_1 = \alpha_2 = \cdots = \alpha_a = 0.$$

若 H_0 被拒绝, 则接受因素 A 的各水平的效应之间有显著差异的假设. 下面来推导 H_0 的检验统计量. 记

$$\bar{y} = \frac{1}{n} \sum_{i=1}^a \sum_{j=1}^{n_i} y_{ij}.$$

考虑统计量

$$SS_T = \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2,$$

称 SS_T 为总离差平方和, 简称为总平方和, 它反映了全部试验数据之间的差异. 对 SS_T 进行分解:

$$\begin{aligned} SS_T &= \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot} + \bar{y}_{i\cdot} - \bar{y})^2 \\ &= \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2 + \sum_{i=1}^a \sum_{j=1}^{n_i} (\bar{y}_{i\cdot} - \bar{y})^2 + 2 \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})(\bar{y}_{i\cdot} - \bar{y}), \end{aligned}$$

其中 $\bar{y}_{i\cdot} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$. 注意到 $\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot}) = 0$, 所以

$$SS_T = \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2 + \sum_{i=1}^a \sum_{j=1}^{n_i} (\bar{y}_{i\cdot} - \bar{y})^2 =: SS_E + SS_A, \quad (7.1.5)$$

其中, $SS_E = \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2$ 反映了随机误差的影响, 因为对固定的 i , $\{y_{i1}, \dots, y_{in_i}\}$ 来自同一个正态总体 $N(\mu_i, \sigma^2)$, 因此它们之间的差异完全是由随机误差所致, 而 $\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2$ 正是它们误差大小的度量. 把 a 组这样的离差平方和求和就得到了 SS_E . 通常称 SS_E 为误差平方和或组内平方和; $SS_A = \sum_{i=1}^a \sum_{j=1}^{n_i} (\bar{y}_{i\cdot} - \bar{y})^2 = \sum_{i=1}^a n_i (\bar{y}_{i\cdot} - \bar{y})^2$. 注意到 $\bar{y}_{i\cdot}$ 是第 i 个总体的样本均值, 它是第 i 个总体均值 μ_i 的估计; \bar{y} 是 $\mu = \frac{1}{n} \sum_{i=1}^a n_i \mu_i$ 的估计. 因此, SS_A 是 a 个总体均值 μ_1, \dots, μ_a 之间的差异大小的一个度量. 称 SS_A 为效应平方和或组间平方和. (7.1.5)是平方和分解公式, 它将总平方和按其来源分解成两部分. 一部分是 SS_E , 即误差平方和, 是由随机误差引起的. 另一部分是 SS_A , 即因素 A 的平方和, 是由因素 A 的各水平的差异引起的.

由于对固定的 i , $\{y_{ij}, j = 1, \dots, n_i\}$ 为来自 $N(\mu_i, \sigma^2)$ 的样本, 因此

$$\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2 / \sigma^2 \sim \chi^2(n_i - 1).$$

所以有

$$E(SS_E) = \sum_{i=1}^a E\left[\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2\right] = (n - a)\sigma^2.$$

这说明 $SS_E/(n-a)$ 是 σ^2 的一个无偏估计. 另一方面,

$$\begin{aligned} E(SS_A) &= E\left[\sum_{i=1}^a n_i (\bar{y}_{i\cdot} - \bar{y} - \alpha_i + \alpha_i)^2\right] \\ &= \sum_{i=1}^a n_i \left[E(\bar{y}_{i\cdot} - \bar{y} - \alpha_i)^2 + \alpha_i^2 \right] \\ &= \sum_{i=1}^a n_i \left(\frac{\sigma^2}{n_i} - \frac{\sigma^2}{n} \right) + \sum_{i=1}^a n_i \alpha_i^2 \\ &= (a-1)\sigma^2 + \sum_{i=1}^a n_i \alpha_i^2. \end{aligned}$$

所以

$$E[SS_A/(a-1)] = \sigma^2 + \sum_{i=1}^a n_i \alpha_i^2 / (a-1).$$

从这个结论可以看出, $SS_A/(a-1)$ 反映了各水平效应的影响. 当 H_0 为真时, $SS_A/(a-1)$ 是 σ^2 的无偏估计. 因此, 若 H_0 为真, 那么

$$F = \frac{SS_A/(a-1)}{SS_E/(n-a)}$$

的取值在1附近; 若 H_0 不真, 则 F 的取值有变大的趋势. 这启发我们可以通过 F 的大小来检验 H_0 .

由样本 $\{y_{ij}, i=1, \dots, a; j=1, \dots, n_i\}$ 的独立性, 可知

$$\frac{SS_E}{\sigma^2} = \frac{\sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2}{\sigma^2} \sim \chi^2(n-a).$$

若 H_0 为真, 那么 $\{y_{ij}, i=1, \dots, a; j=1, \dots, n_i\}$ 是独立同分布随机变量序列, 服从 $N(\mu, \sigma^2)$. 所以 H_0 为真时,

$$\frac{SS_T}{\sigma^2} \sim \chi^2(n-1).$$

为了推导出检验统计量在 H_0 为真时的分布, 我们把 SS_T, SS_E, SS_A 都写成正态随机向量的二次型的形式. 回忆

$$\mathbf{Y} = (y_{11}, \dots, y_{1n_1}, \dots, y_{a1}, \dots, y_{an_a})'.$$

记 $\mathbf{1}_n$ 为元素全为1的 n 维列向量, $\mathbf{1}_{n_i}$ 为元素全为1的 n_i 维列向量, $i=1, \dots, a$. 那么, 若 H_0 为真, 则

$$\mathbf{Y} \sim N(\mu \mathbf{1}_n, \sigma^2 \mathbf{I}_n), \quad \frac{\mathbf{Y}}{\sigma} \sim N\left(\frac{\mu}{\sigma} \mathbf{1}_n, \mathbf{I}_n\right).$$

注意到

$$\begin{cases} SS_T = \mathbf{Y}'(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n') \mathbf{Y} =: \mathbf{Y}' \mathbf{C} \mathbf{Y}, \\ SS_E = \mathbf{Y}'(\mathbf{I}_n - \text{diag}(\frac{1}{n_1} \mathbf{1}_{n_1} \mathbf{1}_{n_1}', \dots, \frac{1}{n_a} \mathbf{1}_{n_a} \mathbf{1}_{n_a}')) \mathbf{Y} =: \mathbf{Y}' \mathbf{C}_1 \mathbf{Y}, \\ SS_A = \mathbf{Y}'(\text{diag}(\frac{1}{n_1} \mathbf{1}_{n_1} \mathbf{1}_{n_1}', \dots, \frac{1}{n_a} \mathbf{1}_{n_a} \mathbf{1}_{n_a}') - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n') \mathbf{Y} =: \mathbf{Y}' \mathbf{C}_2 \mathbf{Y}. \end{cases}$$

现已知:

- (1) $\mathbf{C} = \mathbf{C}_1 + \mathbf{C}_2$;
- (2) \mathbf{C}_2 为非负定矩阵(因为它是对称幂等矩阵, 特征根非0即1);
- (3) $SS_E/\sigma^2 \sim \chi^2(n-a)$;
- (4) 若 H_0 为真, 则 $SS_T/\sigma^2 \sim \chi^2(n-1)$.

那么, 根据定理2.4.4, H_0 为真时, 有

$$\frac{SS_A}{\sigma^2} \sim \chi^2(a-1, \lambda)$$

且 SS_A 与 SS_E 相互独立, 其中非中心参数

$$\lambda = \left(\frac{\mu}{\sigma} \mathbf{1}_n\right)' \mathbf{C}_2 \left(\frac{\mu}{\sigma} \mathbf{1}_n\right) = 0.$$

此外可知, H_0 为真时, 检验统计量

$$F = \frac{SS_A/(a-1)}{SS_E/(n-a)} \sim F(a-1, n-a). \quad (7.1.6)$$

因此, 给定显著性水平 α , 假设检验的拒绝域为

$$\{\text{样本}: F > F_\alpha(a-1, n-a)\}.$$

下面的方差分析表(表7.1.2)描述了上述的假设检验过程.

方差来源	平方和	自由度	均方	F 比
因素A	SS_A	$a-1$	$MS_A = \frac{SS_A}{a-1}$	$F = \frac{MS_A}{MS_E}$
误差	SS_E	$n-a$	$MS_E = \frac{SS_E}{n-a}$	
总和	SS_T	$n-1$		

表7.1.2 方差分析表

例7.1.1 设有三个小麦品种, 经试种得每公顷产量数据(单位: kg/hm^2)见表7.1.3. 问: 不同品种的小麦产量之间有无显著的差异?

品种\试验号	1	2	3	4	5
1	4350	4650	4080	4275	
2	4125	3720	3810	3960	3930
3	4695	4245	4620		

表7.1.3 小麦品种试验数据

计算与分析过程: $\bar{y} = 4205$; $\bar{y}_{1\cdot} = 4339$, $n_1 = 4$; $\bar{y}_{2\cdot} = 3909$, $n_2 = 5$; $\bar{y}_{3\cdot} = 4520$, $n_3 = 3$; 总平方和 $SS_T = 1186800$, 自由度为 $n-1 = 12-1 = 11$; 效

应平方和 $SS_A = 807311.25$, 自由度为 $a - 1 = 2$; 误差平方和 $SS_E = 379488.75$, 自由度为 $n - a = 12 - 3 = 9$. 因此, F 检验统计量的样本值为

$$F = \frac{807311.25/2}{379488.75/9} = 9.57.$$

下面的方差分析表描述了整个计算过程. 给定显著性水平0.05, 因 $F_{0.05}(2, 9) = 4.26 < 9.57$, 所以拒绝原假设, 认为小麦品种的效应具有显著的差异.

方差来源	平方和	自由度	均方	F 比
因素A	807311.25	2	403655.62	9.57
误差	379488.75	9	42165.42	
总和	1186800	11		

表7.1.4 小麦品种的方差分析表

对该例子, R代码及分析结果如下:

```

1 > wheat=data.frame(
2 > X=c(4350,4650,4080,4275,4125,3720,3810,3960,3930,4695,4245,4620),
3 > A=factor(rep(1:3,c(4,5,3)))
4 > )
5 > wheat.aov=aov(X~A,data=wheat)
6 > summary(wheat.aov)
7           Df Sum Sq Mean Sq F value    Pr(>F)
8 A           2  807311   403656    9.573 0.00591 **
9 Residuals    9  379489    42165
10 ---
11 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

由于 p 值= 0.00591 < 0.05, 所以拒绝原假设.

在进行方差分析之前, 对方差的齐性假设进行检验. 可使用Bartlett检验方法, R代码及分析结果如下:

```

1 > bartlett.test(X~A,data=wheat)
2       Bartlett test of homogeneity of variances
3 data:  X by A
4 Bartlett's K-squared = 0.68263, df = 2, p-value = 0.7108

```

由于 p 值= 0.7108 > 0.05, 所以接受方差齐性假设.

如果 F 检验的结论是拒绝原假设, 则表明从现在掌握的数据看, 我们有理由认为因素A的 a 个水平效应之间有显著的差异, 也就是说, μ_1, \dots, μ_a 不完全相同. 这时, 需要对每一对 μ_i 和 μ_j 之间的差异程度作出估计. 这就要对效应之差 $\mu_i - \mu_j$ 作区间估计, 或者对 $H_0: \mu_i = \mu_j$ 进行假设检验. 不难看出:

$$\frac{(\bar{y}_{i\cdot} - \bar{y}_{j\cdot}) - (\mu_i - \mu_j)}{\sigma \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} \sim N(0, 1); \quad (7.1.7)$$

记 $\hat{\sigma}^2 = SS_E/(n-a)$, 那么

$$\frac{(n-a)\hat{\sigma}^2}{\sigma^2} = \frac{SS_E}{\sigma^2} \sim \chi^2(n-a); \quad (7.1.8)$$

且

$(\bar{y}_{i\cdot} - \bar{y}_{j\cdot})$ 与 SS_E 相互独立.

所以

$$\frac{(\bar{y}_{i\cdot} - \bar{y}_{j\cdot}) - (\mu_i - \mu_j)}{\hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} \sim t(n-a).$$

因此在 $H_0: \mu_i = \mu_j$ 成立时, 检验统计量

$$t_{ij} = \frac{\bar{y}_{i\cdot} - \bar{y}_{j\cdot}}{\hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} \sim t(n-a).$$

给定显著性水平 α , 检验的拒绝域为

$$W = \{\text{样本}: |t_{ij}| > t_{\alpha/2}(n-a)\}.$$

或者用区间估计的方法进行统计推断. $\mu_i - \mu_j$ 的置信水平为 $1 - \alpha$ 的置信区间为

$$\left(\bar{y}_{i\cdot} - \bar{y}_{j\cdot} - \hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}} t_{\alpha/2}(n-a), \bar{y}_{i\cdot} - \bar{y}_{j\cdot} + \hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}} t_{\alpha/2}(n-a) \right).$$

如果这个区间包含0, 则表明可以概率 $1 - \alpha$ 断言 μ_i 与 μ_j 没有显著差异; 如果整个区间落在0的左边, 则可以概率 $1 - \alpha$ 断言 μ_i 小于 μ_j ; 如果整个区间落在0的右边, 则可以概率 $1 - \alpha$ 断言 μ_i 大于 μ_j .

例7.1.1续 考察 μ_1, μ_2, μ_3 之间的差异.

先画箱线图进行直观的判断, R代码如下:

```
1 > plot(wheat$X~wheat$A)
```

得到图7.1.1, 可以看出, μ_2 与 μ_3 之间的差异比较明显.

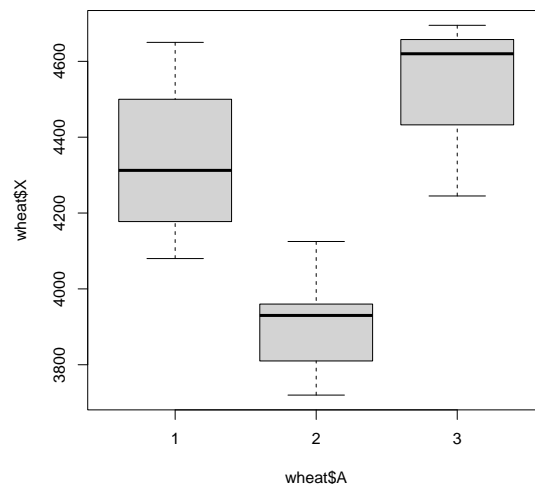


图7.1.1 箱线图

下面先计算各个因子的均值, 然后用多重 t 检验的方法进行统计推断, R代码及分析结果如下:

```

1 > mu=tapply(wheat$X,wheat$A,mean) #tapply命令用于分组统计
2 > mu
3      1      2      3
4 4338.75 3909.00 4520.00
5 > pairwise.t.test(wheat$X,wheat$A,p.adjust.method="none")
6 Pairwise comparisons using t tests with pooled SD
7 data:  wheat$X and wheat$A
8      1      2
9 2 0.0123 -
10 3 0.2776 0.0028
11 P value adjustment method: none

```

根据 p 值, 我们以95%的把握断言: μ_1 和 μ_2 有显著差异, μ_2 和 μ_3 有显著差异, μ_1 和 μ_3 无显著差异.

但要注意的是, 现在进行的是多个假设检验, 这些 p 值是不可靠的, 为什么呢? 假设 μ_1, μ_2, μ_3 之间并无显著差异, 那么两两比较的显著性检验的 p 值在区间 $[0, 1]$ 上呈现均匀分布. 因此, 至少出现一个虚假的显著性结果的概率为

$$1 - 0.95^3 = 0.142.$$

若同时进行6个两两比较的显著性检验, 则至少出现一个虚假的显著性结果的概率上升为

$$1 - 0.95^6 = 0.265.$$

所以有必要对 p 值进行修正. 其中的一个修正方法是Bonferroni修正法. 假设需要进行两两比较的假设检验个数是 m 个, 那么Bonferroni修正法把满足 $m \times p\text{值} < \alpha$ 的结果才称为是统计显著的(即有显著性差异). 显然, Bonferroni修正法过于保守了. 另一种常用的修正方法是FDR(False Discovery Rate)方法, 它由Benjamini和Hochberg在1995年提出. 该方法先对 m 个 p 值进行排序:

$$p_{(1)} \leq \cdots \leq p_{(m)},$$

然后选择满足

$$p_{(i)} \leq \alpha i / m$$

的最大 i 值, 并把 $p_{(1)}, \dots, p_{(i)}$ 所对应的结果称为是统计显著的(即有显著性差异).

再回到前面的例子, 我们对两两 t 检验的 p 值进行修正. R代码及分析结果如下:

```

1 > pairwise.t.test(wheat$X,wheat$A,p.adjust.method="bonferroni")
2   Pairwise comparisons using t tests with pooled SD
3 data:  wheat$X and wheat$A
4    1      2
5 2 0.0370 -
6 3 0.8327 0.0083
7 P value adjustment method: bonferroni
8 > pairwise.t.test(wheat$X,wheat$A,p.adjust.method="fdr")
9   Pairwise comparisons using t tests with pooled SD
10 data:  wheat$X and wheat$A
11    1      2
12 2 0.0185 -
13 3 0.2776 0.0083
14 P value adjustment method: fdr

```

可以看出, 若取 $\alpha = 0.05$, Bonferroni方法和FDR方法仍认为 μ_1 和 μ_2 有显著差异, μ_2 和 μ_3 有显著差异. 但若取 $\alpha = 0.03$, 则Bonferroni方法只认为 μ_2 和 μ_3 有显著差异, 而FDR方法仍维持 $\alpha = 0.05$ 时的判断.

7.2 两因素方差分析

在一项实际试验中, 往往有这样的情况, 研究者本想考察某个因素对指标的影响, 但是由于客观条件的限制, 还有个别因素不可能在所有试验中把它们控制在完全相同的状态. 譬如, 在上一节开始讨论的小麦品种的例子中, 实验者要研究的是“小麦品种”这一因素对产量的影响. 但在实际中可能会出现这样的情况, 很难找到一大块田, 其土质肥沃程度完全一样. 因此“土质”就成为另一个因素不可避免地进入了试验, 导致了两因素的试验问题.

在农业试验中解决这个问题的方法是采用所谓的区组设计. 它的做法是, 先把一块田(大块)分成若干块(小块), 譬如 b (小块), 使得同一大块田里的各小块的土质肥沃程度基本上保持一致. 在试验设计中, 称这种大块(田)为区组, 然后把每一个区组又分成若干小块(田), 称为试验单元. 现在有 a 个小麦品种, 一个方便的做法就是把每个区组(大块)分成 a 个试验单元(小块). 在每一个试验单元上种

植一种小麦. 若用 y_{ij} 表示在第 j 个区组中种植第 i 种小麦的那个试验单元的产量, 则 y_{ij} 就可表为

$$y_{ij} = \mu + \alpha_i + \beta_j + e_{ij}, \quad i = 1, \dots, a; \quad j = 1, \dots, b, \quad (7.2.1)$$

这里, μ 称为总平均, α_i 为第 i 个小麦品种的效应, β_j 为第 j 个区组的效应, e_{ij} 为随机误差.

在实际应用中, 更多的情况是研究者所感兴趣的问题本身就是两因素的. 例如, 在一项工业试验中, 影响产品质量的因素是反应温度和反应压力. 试验者的目的是选择最好的生产条件, 若反应温度有 a 个水平, 反应压力有 b 个水平, 记 y_{ij} 为在反应温度处于第 i 个水平和反应压力处于第 j 个水平时产品质量的指标值, 那么 y_{ij} 也有表达式(7.2.1). 若影响产品质量的因素除了反应温度和反应压力外, 还有反应时间和催化剂这两个因素, 当我们把任何两个因素控制在某一状态而研究剩余两个因素对产品质量的影响时, 这同样导致一个两因素的试验问题.

考虑一般的两因素试验问题, 将这两个因素分别记为 A 和 B . 假定因素 A 有 a 个不同的水平, 记为 A_1, \dots, A_a , 因素 B 有 b 个不同的水平, 记为 B_1, \dots, B_b . 在因素 A 和 B 的各个水平的组合下做 c 次试验. 记 y_{ijk} 为在水平组合 (A_i, B_j) 下的第 k 次试验的指标值. 对固定的 i 和 j , $\{y_{ij1}, \dots, y_{ijc}\}$ 都是在水平组合 (A_i, B_j) 下的指标观测值, 因此可以把它们看成来自同一个正态总体的样本, 均值记为 μ_{ij} . 于是

$$\text{给定 } i, j, \quad y_{ijk} \stackrel{i.i.d.}{\sim} N(\mu_{ij}, \sigma^2), \quad k = 1, \dots, c. \quad (7.2.2)$$

表7.2.1描述了两因素方差分析问题的总体和样本.

$A_i \setminus B_j$	B_1	B_2	\dots	B_b
A_1	$y_{111}, y_{112}, \dots, y_{11c}$	$y_{121}, y_{122}, \dots, y_{12c}$	\dots	$y_{1b1}, y_{1b2}, \dots, y_{1bc}$
A_2	$y_{211}, y_{212}, \dots, y_{21c}$	$y_{221}, y_{222}, \dots, y_{22c}$	\dots	$y_{2b1}, y_{2b2}, \dots, y_{2bc}$
\vdots	\vdots	\vdots	\ddots	\vdots
A_a	$y_{a11}, y_{a12}, \dots, y_{a1c}$	$y_{a21}, y_{a22}, \dots, y_{a2c}$	\dots	$y_{ab1}, y_{ab2}, \dots, y_{abc}$

表7.2.1 两因素方差分析问题

将(7.2.2)改写成

$$\begin{cases} y_{ijk} = \mu_{ij} + e_{ijk}, & i = 1, \dots, a; \quad j = 1, \dots, b; \quad k = 1, \dots, c, \\ e_{ijk} \stackrel{i.i.d.}{\sim} N(0, \sigma^2). \end{cases} \quad (7.2.3)$$

为进行统计分析, 将 μ_{ij} 做适当的分解. 记

$$\begin{aligned} \mu &= \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b \mu_{ij}, \quad \bar{\mu}_{i\cdot} = \frac{1}{b} \sum_{j=1}^b \mu_{ij}, \quad \bar{\mu}_{\cdot j} = \frac{1}{a} \sum_{i=1}^a \mu_{ij}, \\ \alpha_i &= \bar{\mu}_{i\cdot} - \mu, \quad i = 1, \dots, a, \\ \beta_j &= \bar{\mu}_{\cdot j} - \mu, \quad j = 1, \dots, b, \\ \gamma_{ij} &= \mu_{ij} - \bar{\mu}_{i\cdot} - \bar{\mu}_{\cdot j} + \mu, \quad i = 1, \dots, a; \quad j = 1, \dots, b, \end{aligned}$$

其中 μ 为总平均, α_i 为因素 A 的水平 A_i 的效应, β_j 为因素 B 的水平 B_j 的效应, γ_{ij} 可写为

$$\gamma_{ij} = \mu_{ij} - (\bar{\mu}_{i\cdot} - \mu) - (\bar{\mu}_{\cdot j} - \mu) - \mu = (\mu_{ij} - \mu) - \alpha_i - \beta_j,$$

表示 A_i 和 B_j 的交互效应. 通常把因素 A 和 B 对试验指标的交互效应设想为某一因素的效应, 称这个因素为 A 与 B 的交互作用, 记为 $A \times B$. 不难验证

$$\sum_{i=1}^a \alpha_i = 0, \quad \sum_{j=1}^b \beta_j = 0, \quad \sum_{i=1}^a \sum_{j=1}^b \gamma_{ij} = 0.$$

注意到 μ_{ij} 可改写为 $\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$, 因此(7.2.3)可写成

$$\begin{cases} y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ijk}, & i = 1, \dots, a; \quad j = 1, \dots, b; \quad k = 1, \dots, c, \\ e_{ijk} \stackrel{i.i.d.}{\sim} N(0, \sigma^2), \\ \sum_{i=1}^a \alpha_i = 0, \quad \sum_{j=1}^b \beta_j = 0, \quad \sum_{i=1}^a \sum_{j=1}^b \gamma_{ij} = 0. \end{cases} \quad (7.2.4)$$

这就是两因素方差分析模型.

注 其实, $\sum_{i=1}^a \sum_{j=1}^b \gamma_{ij} = 0$ 应写成

$$\sum_{j=1}^b \gamma_{ij} = 0, \quad \forall i = 1, \dots, a; \quad \sum_{i=1}^a \gamma_{ij} = 0, \quad \forall j = 1, \dots, b.$$

这里共有 $a + b - 1$ 个线性约束.

7.2.1 无交互效应的情形

假设对所有的 $i = 1, \dots, a$ 和 $j = 1, \dots, b$, 都有 $\gamma_{ij} = 0$, 即不存在交互效应. 此时, 只考虑每种水平组合下试验次数为 $c = 1$ 的情形. 模型(7.2.4)可写为

$$\begin{cases} y_{ij} = \mu + \alpha_i + \beta_j + e_{ij}, & i = 1, \dots, a; \quad j = 1, \dots, b, \\ e_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2), \\ \sum_{i=1}^a \alpha_i = 0, \quad \sum_{j=1}^b \beta_j = 0. \end{cases} \quad (7.2.5)$$

这就是无交互效应的两因素方差分析模型. 试验的目的是考察因素 A 或 B 的各个水平对试验指标的影响有无显著差异, 这归结为对假设

$$H_1 : \alpha_1 = \alpha_2 = \dots = \alpha_a = 0$$

和

$$H_2 : \beta_1 = \beta_2 = \dots = \beta_b = 0$$

进行检验. 我们采用与单因素方差分析模型类似的方法导出检验统计量.

记

$$\bar{y} = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b y_{ij}, \quad \bar{y}_{i\cdot} = \frac{1}{b} \sum_{j=1}^b y_{ij}, \quad \bar{y}_{\cdot j} = \frac{1}{a} \sum_{i=1}^a y_{ij},$$

$$SS_T = \sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y})^2.$$

SS_T 是全部试验数据的离差平方和, 称为总平方和. 对其进行分解得

$$\begin{aligned} SS_T &= \sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{\cdot j} + \bar{y} + \bar{y}_{i\cdot} - \bar{y} + \bar{y}_{\cdot j} - \bar{y})^2 \\ &= \sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{\cdot j} + \bar{y})^2 + \sum_{i=1}^a b(\bar{y}_{i\cdot} - \bar{y})^2 + \sum_{j=1}^b a(\bar{y}_{\cdot j} - \bar{y})^2 \\ &=: SS_E + SS_A + SS_B. \end{aligned}$$

因为 $\bar{y}_{i\cdot}$ 是水平 A_i 下所有观测值的平均, 所以 $\sum_{i=1}^a (\bar{y}_{i\cdot} - \bar{y})^2$ 反映了 $\bar{y}_{1\cdot}, \dots, \bar{y}_{a\cdot}$ 之间差异的程度. 这种差异是由于因素 A 的不同水平所引起的, 因此称 SS_A 为因素 A 的平方和. 类似地, 称 SS_B 为因素 B 的平方和. SS_E 反映了试验的随机误差的影响, 称为误差平方和. 事实上, 可写

$$\begin{aligned} SS_E &= \sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{\cdot j} + \bar{y})^2 \\ &= \sum_{i=1}^a \sum_{j=1}^b [y_{ij} - (\bar{y}_{i\cdot} - \bar{y}) - (\bar{y}_{\cdot j} - \bar{y}) - \bar{y}]^2, \end{aligned}$$

而 $(\bar{y}_{i\cdot} - \bar{y})$ 的均值为 $\bar{\mu}_{i\cdot} - \mu = \alpha_i$, $(\bar{y}_{\cdot j} - \bar{y})$ 的均值为 $\bar{\mu}_{\cdot j} - \mu = \beta_j$, \bar{y} 的均值为 μ , 所以可把 $y_{ij} - (\bar{y}_{i\cdot} - \bar{y}) - (\bar{y}_{\cdot j} - \bar{y}) - \bar{y}$ 近似看作 $y_{ij} - \mu - \alpha_i - \beta_j = e_{ij}$.

可以证明:

$$SS_E/\sigma^2 \sim \chi^2((a-1)(b-1));$$

当 H_1 成立时,

$$SS_A/\sigma^2 \sim \chi^2(a-1) \text{ 且与 } SS_E \text{ 独立};$$

当 H_2 成立时,

$$SS_B/\sigma^2 \sim \chi^2(b-1) \text{ 且与 } SS_E \text{ 独立}.$$

注 (1) 证明 $SS_E/\sigma^2 \sim \chi^2((a-1)(b-1))$ 的关键是把 SS_E 写成

$$\begin{aligned} SS_E &= \sum_{i=1}^a \sum_{j=1}^b (e_{ij} - \bar{e}_{i\cdot} - \bar{e}_{\cdot j} + \bar{e})^2 \\ &= \sum_{i=1}^a \sum_{j=1}^b (e_{ij} - \bar{e}_{i\cdot})^2 - \sum_{i=1}^a \sum_{j=1}^b (\bar{e}_{\cdot j} - \bar{e})^2, \end{aligned}$$

其中

$$\bar{e} = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b e_{ij}, \quad \bar{e}_{i\cdot} = \frac{1}{b} \sum_{j=1}^b e_{ij}, \quad \bar{e}_{\cdot j} = \frac{1}{a} \sum_{i=1}^a e_{ij}.$$

因此 $\sum_{i=1}^a \sum_{j=1}^b (e_{ij} - \bar{e}_{i.})^2$ 有如下的平方和分解式:

$$\sum_{i=1}^a \sum_{j=1}^b (e_{ij} - \bar{e}_{i.})^2 = \sum_{i=1}^a \sum_{j=1}^b (e_{ij} - \bar{e}_{i.} - \bar{e}_{.j} + \bar{e})^2 + \sum_{i=1}^a \sum_{j=1}^b (\bar{e}_{.j} - \bar{e})^2.$$

(2) 证明当 H_1 成立时, $SS_A/\sigma^2 \sim \chi^2(a-1)$ 的关键是写 SS_A 为

$$SS_A = \sum_{i=1}^a b(\bar{y}_{i.} - \bar{y})^2 = \sum_{i=1}^a (\sqrt{b}\bar{y}_{i.} - \sqrt{b}\bar{y})^2,$$

其中 $\sqrt{b}\bar{y}_{i.} \sim N(\mu + \alpha_i, \sigma^2)$. 所以当 H_1 成立时, $\sqrt{b}\bar{y}_{i.} \sim N(\mu, \sigma^2)$, $\{\sqrt{b}\bar{y}_{i.}, i = 1, \dots, a\}$ 为来自 $N(\mu, \sigma^2)$ 的样本.

(3) 证明当 H_2 成立时, $SS_B/\sigma^2 \sim \chi^2(b-1)$ 的关键是写 SS_B 为

$$SS_B = \sum_{j=1}^b a(\bar{y}_{.j} - \bar{y})^2 = \sum_{j=1}^b (\sqrt{a}\bar{y}_{.j} - \sqrt{a}\bar{y})^2,$$

其中 $\sqrt{a}\bar{y}_{.j} \sim N(\mu + \beta_j, \sigma^2)$. 所以当 H_2 成立时, $\sqrt{a}\bar{y}_{.j} \sim N(\mu, \sigma^2)$, $\{\sqrt{a}\bar{y}_{.j}, j = 1, \dots, b\}$ 为来自 $N(\mu, \sigma^2)$ 的样本.

于是当 H_1 成立时,

$$F_A = \frac{SS_A/(a-1)}{SS_E/[(a-1)(b-1)]} \sim F(a-1, (a-1)(b-1)). \quad (7.2.6)$$

给定显著性水平 α , 假设检验的拒绝域为

$$W = \{\text{样本} : F_A > F_{\alpha}(a-1, (a-1)(b-1))\}.$$

同理, 当 H_2 成立时,

$$F_B = \frac{SS_B/(b-1)}{SS_E/[(a-1)(b-1)]} \sim F(b-1, (a-1)(b-1)). \quad (7.2.7)$$

给定显著性水平 α , 假设检验的拒绝域为

$$W = \{\text{样本} : F_B > F_{\alpha}(b-1, (a-1)(b-1))\}.$$

下面的两因素方差分析表(表7.2.2)描述了上述的分析过程.

方差来源	平方和	自由度	均方	F比
因素A	SS_A	$a-1$	$MS_A = \frac{SS_A}{a-1}$	$F_A = \frac{MS_A}{MS_E}$
因素B	SS_B	$b-1$	$MS_B = \frac{SS_B}{b-1}$	$F_B = \frac{MS_B}{MS_E}$
误差	SS_E	$(a-1)(b-1)$	$MS_E = \frac{SS_E}{(a-1)(b-1)}$	
总计	SS_T	$ab-1$		

表7.2.2 无交互效应的两因素方差分析表

例7.2.1 一种火箭使用了四种燃料、三种推进器, 进行射程试验. 对于每种燃料和推进器的组合做一次试验, 得到试验数据见表7.2.3. 问在0.05的显著性水平下各种燃料之间及各种推进器之间有无显著差异?

燃料A\推进器B	B_1	B_2	B_3
A_1	58.2	56.2	65.3
A_2	49.1	54.1	51.6
A_3	60.1	70.9	39.23
A_4	75.8	58.2	48.7

表7.2.3 火箭试验数据

这是一个两因素试验且不考虑交互效应. 记燃料为因素A, 它有4个水平, 水平效应为 $\alpha_1, \dots, \alpha_4$; 推进器为因素B, 它有3个水平, 水平效应为 $\beta_1, \beta_2, \beta_3$. 我们在显著性水平 $\alpha = 0.05$ 下检验

$$H_1: \alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 0,$$

$$H_2: \beta_1 = \beta_2 = \beta_3 = 0.$$

经计算(略)得下面的方差分析表:

方差来源	平方和	自由度	均方	F比
因素A	157.59	3	52.53	$F_A = 0.4306$
因素B	223.85	2	111.93	$F_B = 0.9175$
误差	731.98	6	122.00	
总计	1113.42	11		

表7.2.4 火箭试验的方差分析表

因为 $F_{0.05}(3, 6) = 4.76$, $F_{0.05}(2, 6) = 5.14$, 而 $F_A = 0.4306 < 4.76$, $F_B = 0.9175 < 5.14$, 所以我们接受 H_1 和 H_2 , 即认为各种燃料和各种推进器之间的差异对于火箭射程无显著影响. 相应的R代码及分析结果如下:

```

1 > rocket=data.frame(
2 > Y=c(58.2,56.2,65.3,49.1,54.1,51.6,60.1,70.9,39.23,75.8,58.2,48.7),
3 > A=gl(4,3),
4 > B=gl(3,1,12)
5 > )
6 > rocket.aov=aov(Y~A+B,data=rocket)
7 > summary(rocket.aov)
8           Df Sum Sq Mean Sq F value Pr(>F)
9 A             3   157.6    52.52   0.431  0.739
10 B             2   223.5   111.74   0.917  0.449
11 Residuals     6   731.3   121.88

```

在上述代码中, gl(4,3)用来产生因子, 该因子有4个水平, 每个水平重复3次. gl(3,1,12)也用来产生因子, 该因子有3个水平, 每个水平重复1次, 循环产生12个因子值. 由于 p 值= 0.739 > 0.05, p 值= 0.449 > 0.05, 所以接受 H_1 和 H_2 .

如果经过 F 检验, H_1 被拒绝, 那么在这种情况下, 我们认为因素 A 的 a 个水平效应 $\alpha_1, \dots, \alpha_a$ 不全相同. 此时我们希望比较 α_i 的大小, 这需要做 $H_0: \alpha_i = \alpha_k$ 的假设检验或者求 $\alpha_i - \alpha_k$ 区间估计. 因为 $y_{ij} \sim N(\mu + \alpha_i + \beta_j, \sigma^2)$, 利用 $\sum_{j=1}^b \beta_j = 0$ 易知

$$\bar{y}_{i\cdot} \sim N(\mu + \alpha_i, \frac{\sigma^2}{b}), \quad i = 1, \dots, a.$$

于是

$$\bar{y}_{i\cdot} - \bar{y}_{k\cdot} \sim N(\alpha_i - \alpha_k, \frac{2\sigma^2}{b}). \quad (7.2.8)$$

注意到

$$\frac{SS_E}{\sigma^2} \sim \chi^2((a-1)(b-1)),$$

且与 $\bar{y}_{i\cdot} - \bar{y}_{k\cdot}$ 相互独立. 因此对固定的 i 和 k , $H_0: \alpha_i = \alpha_k$ 的检验统计量

$$t_{ik} = \frac{\sqrt{b}(\bar{y}_{i\cdot} - \bar{y}_{k\cdot})}{\sqrt{2}\hat{\sigma}} \stackrel{H_0}{\sim} t((a-1)(b-1)).$$

给定显著性水平 α , 检验的拒绝域为

$$W = \{\text{样本} : |t_{ik}| > t_{\alpha/2}((a-1)(b-1))\}.$$

或者考虑区间估计, $\alpha_i - \alpha_k$ 的置信水平为 $1 - \alpha$ 的置信区间为

$$\left(\bar{y}_{i\cdot} - \bar{y}_{k\cdot} \pm \sqrt{\frac{2}{b}} \hat{\sigma} t_{\alpha/2}((a-1)(b-1)) \right). \quad (7.2.9)$$

如果这个区间包含0, 则表明我们可以概率 $1 - \alpha$ 断言 α_i 和 α_k 没有显著差异. 如果整个区间落在0的左边, 则可以概率 $1 - \alpha$ 断言 α_i 小于 α_k . 如果整个区间落在0的右边, 则可以概率 $1 - \alpha$ 断言 α_i 大于 α_k .

若经过 F 检验, 假设 H_2 被拒绝, 类似于以上的讨论, 可以构建 $H_0: \beta_j = \beta_k$ 的检验统计量

$$t_{jk} = \frac{\sqrt{a}(\bar{y}_{\cdot j} - \bar{y}_{\cdot k})}{\sqrt{2}\hat{\sigma}} \stackrel{H_0}{\sim} t((a-1)(b-1)).$$

给定显著性水平 α , 检验的拒绝域为

$$W = \{\text{样本} : |t_{jk}| > t_{\frac{\alpha}{2}}((a-1)(b-1))\}.$$

或者构建 $\beta_j - \beta_k$ 的置信水平为 $1 - \alpha$ 的置信区间:

$$\left(\bar{y}_{\cdot j} - \bar{y}_{\cdot k} \pm \sqrt{\frac{2}{a}} \hat{\sigma} t_{\frac{\alpha}{2}}((a-1)(b-1)) \right). \quad (7.2.10)$$

如果这个区间包含0, 则表明可以概率 $1 - \alpha$ 断言 β_j 和 β_k 没有显著差异. 如果整个区间落在0的左边, 则可以概率 $1 - \alpha$ 断言 β_j 小于 β_k . 如果整个区间落在0的右边, 则可以概率 $1 - \alpha$ 断言 β_j 大于 β_k .

与单因素方差分析类似, 做多重假设检验的时候, 需要对 p 值进行调整. 对刚才的例子, 若 H_1 被拒绝(其实并没有被拒绝), 则还需进行多重假设检验, 分析哪些水平之间有显著差异. R代码及分析结果如下:

```
1 > pairwise.t.test(rocket$Y,rocket$A,p.adjust.method="fdr")
2 Pairwise comparisons using t tests with pooled SD
3 data: rocket$Y and rocket$A
4 1 2 3
5 2 0.88 - -
6 3 0.88 0.88 -
7 4 0.91 0.88 0.88
8 P value adjustment method: fdr
```

同样地, 若 H_2 被拒绝(其实并没有被拒绝), 也还需进行多重假设检验, 分析哪些水平之间有显著差异. R代码及分析结果如下:

```
1 > pairwise.t.test(rocket$Y,rocket$B,p.adjust.method="fdr")
2 Pairwise comparisons using t tests with pooled SD
3 data: rocket$Y and rocket$B
4 1 2
5 2 0.90 -
6 3 0.37 0.37
7 P value adjustment method: fdr
```

7.2.2 有交互效应的情形

若要考虑因素 A, B 之间的交互作用 $A \times B$, 那么在各水平组合下需要做重复试验. 假设每种组合下的试验次数均为 $c(c > 1)$. 此时对应的统计模型就是(7.2.4). 在这样的模型下, α_i 并不能反映水平 A_i 的优劣. 这是因为在交互效应存在的情况下, 因子水平 A_i 的优劣还与因子 B 的水平有关系. 因此, 对这样的模型, 直接检验 $\alpha_1 = \dots = \alpha_a = 0$ 与 $\beta_1 = \dots = \beta_b = 0$ 都是没有实际意义的. 一个重要的检验问题是交互效应是否存在, 即检验

$$H_3: \gamma_{ij} = 0, \quad i = 1, \dots, a; \quad j = 1, \dots, b.$$

若 H_3 被接受, 那么检验 $\alpha_1 = \dots = \alpha_a = 0$ 与检验 $\beta_1 = \dots = \beta_b = 0$ 才有意义.

引进记号:

$$\begin{aligned} \bar{y} &= \frac{1}{abc} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c y_{ijk}, & \bar{y}_{ij\cdot} &= \frac{1}{c} \sum_{k=1}^c y_{ijk}, \\ \bar{y}_{i\cdot\cdot} &= \frac{1}{bc} \sum_{j=1}^b \sum_{k=1}^c y_{ijk}, & \bar{y}_{\cdot j\cdot} &= \frac{1}{ac} \sum_{i=1}^a \sum_{k=1}^c y_{ijk}. \end{aligned}$$

作平方和分解:

$$\begin{aligned}
SS_T &= \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (y_{ijk} - \bar{y})^2 \\
&= \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (y_{ijk} - \bar{y}_{ij\cdot} + \bar{y}_{i\cdot\cdot} - \bar{y} + \bar{y}_{\cdot j\cdot} - \bar{y} + \bar{y}_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y})^2 \\
&= \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (y_{ijk} - \bar{y}_{ij\cdot})^2 + bc \sum_{i=1}^a (\bar{y}_{i\cdot\cdot} - \bar{y})^2 + ac \sum_{j=1}^b (\bar{y}_{\cdot j\cdot} - \bar{y})^2 \\
&\quad + c \sum_{i=1}^a \sum_{j=1}^b (\bar{y}_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y})^2 \\
&=: SS_E + SS_A + SS_B + SS_{A \times B},
\end{aligned}$$

其中

$$\begin{aligned}
SS_E &= \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (y_{ijk} - \bar{y}_{ij\cdot})^2, \\
SS_A &= bc \sum_{i=1}^a (\bar{y}_{i\cdot\cdot} - \bar{y})^2, \\
SS_B &= ac \sum_{j=1}^b (\bar{y}_{\cdot j\cdot} - \bar{y})^2, \\
SS_{A \times B} &= c \sum_{i=1}^a \sum_{j=1}^b (\bar{y}_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y})^2.
\end{aligned}$$

称 SS_E 为误差平方和, SS_A 为因素 A 的平方和, SS_B 为因素 B 的平方和, $SS_{A \times B}$ 为交互作用的平方和. 类似于以前的讨论, 可以证明当 H_3 成立时,

$$F_{A \times B} = \frac{SS_{A \times B} / [(a-1)(b-1)]}{SS_E / [ab(c-1)]} \sim F((a-1)(b-1), ab(c-1)). \quad (7.2.11)$$

因此, 给定显著性水平 α , 假设检验的拒绝域为

$$W = \{ \text{样本} : F_{A \times B} > F_\alpha((a-1)(b-1), ab(c-1)) \}.$$

下面的方差分析表(表7.2.5)描述了上述分析过程.

方差来源	平方和	自由度	均方	F比
因素A	SS_A	$a - 1$	$MS_A = \frac{SS_A}{a-1}$	$F_A = \frac{MS_A}{MS_E}$
因素B	SS_B	$b - 1$	$MS_B = \frac{SS_B}{b-1}$	$F_B = \frac{MS_B}{MS_E}$
交互效应A × B	$SS_{A \times B}$	$(a - 1)(b - 1)$	$MS_{A \times B} = \frac{SS_{A \times B}}{(a-1)(b-1)}$	$F_{A \times B} = \frac{MS_{A \times B}}{MS_E}$
误差	SS_E	$ab(c - 1)$	$MS_E = \frac{SS_E}{ab(c-1)}$	
总计	SS_T	$abc - 1$		

表7.2.5 关于交互效应的两因素方差分析表

例7.2.2 研究树种与地理位置对松树生长的影响, 对4个地区的3种同龄松树的直径进行测量得到数据, 见表7.2.6. A_1, A_2, A_3 表示三个不同的树种, B_1, B_2, B_3, B_4 表示4个不同的地区. 对每一种水平组合, 进行了5次测量, 对比试验结果进行方差分析.

	B_1	B_2	B_3	B_4
A_1	23,25,21,14,15	20,17,11,26,21	16,19,13,16,24	20,21,18,27,24
A_2	28,30,19,17,22	26,24,21,25,26	19,18,19,20,25	26,26,28,29,23
A_3	18,15,23,18,10	21,25,12,12,22	19,23,22,14,13	22,13,12,22,19

表7.2.6 3种同龄松树的直径测量数据(单位: cm)

R代码及分析结果如下:

```

1 > tree<-data.frame(
2 +   A=gl(3,20,60),
3 +   B=gl(4,5,60),
4 +   Y=c(23, 25, 21, 14, 15, 20, 17, 11, 26, 21,
5 +       16, 19, 13, 16, 24, 20, 21, 18, 27, 24,
6 +       28, 30, 19, 17, 22, 26, 24, 21, 25, 26,
7 +       19, 18, 19, 20, 25, 26, 26, 28, 29, 23,
8 +       18, 15, 23, 18, 10, 21, 25, 12, 12, 22,
9 +       19, 23, 22, 14, 13, 22, 13, 12, 22, 19)
10 + )
11 > tree.aov=aov(Y~A+B+A:B, data=tree)
12 > summary(tree.aov)
13           Df Sum Sq Mean Sq F value    Pr(>F)
14 A           2   352.5    176.27   8.959 0.000494 ***
15 B           3    87.5     29.17   1.483 0.231077
16 A:B         6    71.7     11.96   0.608 0.722890
17 Residuals  48   944.4     19.68
18 ---
19 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

可见在显著性水平 $\alpha = 0.05$ 下, 树种和地理位置的交互效应不显著(因为 p 值=0.722890 > 0.05). 接下来分别考察树种的效应和地理位置的效应. 由分析结果看出, 树种(因素 A)的效应是显著的(因为 p 值= 0.000494 < 0.05), 而地理位置(因素 B)的效应并不显著(因为 p 值= 0.231077 > 0.05). 最后来看一下树种的多重假设检验(跟前面一样, 需要对 p 值进行调整).

```
1 > pairwise.t.test(tree$Y,tree$A,p.adjust.method="fdr")
2 Pairwise comparisons using t tests with pooled SD
3 data: tree$Y and tree$A
4 1 2
5 2 0.00852 -
6 3 0.20103 0.00032
7 P value adjustment method: fdr
```

由 p 值可以看出, 在显著性水平 $\alpha = 0.05$ 下, A_1 和 A_2 的效应有显著差异, A_2 和 A_3 的效应有显著差异, 而 A_1 和 A_3 的效应无显著差异.

作业

1. 为了比较消费者对于四种打车软件的满意程度, 随机抽取了32人, 分为四组, 让他们分别使用这四种打车软件并给出满意度打分(0 – 100分), 具体数据见下表.

(1) 将样本表示成单因素方差分析模型的形式;

(2) 在显著性水平 $\alpha = 0.05$ 下检验人们对于这四种打车软件的满意程度是否存在显著差异.

打车软件	满意度评分							
第一种	60	91	84	58	91	78	75	63
第二种	68	66	62	65	72	79	70	72
第三种	78	86	84	82	100	97	89	78
第四种	72	66	75	56	58	73	71	66

题1的数据

2. 某乳制品企业有三个车间加工生产250毫升的盒装牛奶. 为了考察三个车间生产的牛奶的蛋白质含量是否一致, 在每个车间生产的产品中各随机抽取了一些样本进行测定, 结果见下表.

(1) 将样本表示成单因素方差分析模型的形式;

(2) 在显著性水平 $\alpha = 0.05$ 下检验这三个车间生产的牛奶的蛋白质含量是否存在显著差异.

车间	蛋白质含量					
车间1	2.41	2.39	2.36	2.43	2.40	2.40
车间2	2.35	2.46	2.39	2.40	2.42	
车间3	2.28	2.37	2.33	2.32	2.24	

题2的数据

3. 有四名工人(W_1, W_2, W_3, W_4)分别操作机床(A_1, A_2, A_3)各一天, 生产同样的产品, 其日产量(单位: 件)见下表. 假设以上两个因素无交互效应.

- (1) 将样本表示成无交互效应的两因素方差分析模型的形式;
- (2) 在显著性水平 $\alpha = 0.05$ 下检验四名工人的日产量有无显著差异;
- (3) 在显著性水平 $\alpha = 0.05$ 下检验三台机床的日产量有无显著差异.

机床\工人	W_1	W_2	W_3	W_4
A_1	50	47	47	53
A_2	63	54	57	58
A_3	52	42	41	48

题3的数据

4. 某医疗机构有五种治疗某癌症的药物治疗方案, 分别记为 A_1, A_2, A_3, A_4, A_5 ; 药物使用剂量有三种方案: 3vg/ml, 10vg/ml和30vg/ml, 分别记为 B_1, B_2, B_3 . 为了分析这两个因素对于癌症治疗的疗效, 在不同的药物治疗方案和剂量组合下, 各进行了1次试验, 得到的癌细胞存活率(%)数据见下表. 假设以上的两个因素没有交互效应.

- (1) 将样本表示成无交互效应的两因素方差分析模型的形式;
- (2) 在显著性水平 $\alpha = 0.05$ 下检验药物治疗方案对癌细胞存活率是否有显著影响;
- (3) 在显著性水平 $\alpha = 0.05$ 下检验药物使用剂量对癌细胞存活率是否有显著影响.

药物\剂量	B_1	B_2	B_3
A_1	90.348	64.933	49.721
A_2	31.028	16.473	27.257
A_3	62.058	41.689	17.343
A_4	17.549	14.902	18.340
A_5	32.334	20.906	17.193

题4的数据

5. 继续讨论题4, 现在在不同的药物和剂量组合下各重复三次试验, 得到的癌细胞存活率(%)数据见下表.

- (1) 把样本表示成有交互效应的两因素方差分析模型的形式;

(2) 在显著性水平 $\alpha = 0.05$ 下检验药物和剂量的交互效应是否显著.

药物\剂量	B1	B2	B3
A1	90.348,81.544,90.954	64.933,63.003,75.093	49.721,46.850,38.023
A2	31.028,37.735,40.991	16.473,31.035,25.001	27.257,12.749,14.509
A3	62.058,68.123,71.225	41.689,39.203,41.074	17.343,22.934,23.456
A4	17.549,24.547,26.994	14.902,17.228,15.773	18.340,11.780,18.411
A5	32.334,34.005,23.907	20.906,25.018,26.992	17.193,16.579,25.210

题5的数据

6. 下表给出了某个化学试验在三种浓度(%)和四种温度($^{\circ}\text{C}$)下的得率(%)的观测值.

- (1) 把样本表示成有交互效应的两因素方差分析模型的形式;
- (2) 在显著性水平 $\alpha = 0.05$ 下检验浓度和温度的交互效应是否显著.
- (3) 在显著性水平 $\alpha = 0.05$ 下检验浓度对产品得率的影响有无显著差异;
- (4) 在显著性水平 $\alpha = 0.05$ 下检验温度对产品得率的影响有无显著差异.

浓度\温度	10	24	38	52
2	14, 10	11, 11	13, 19	10, 12
4	9, 7	10, 8	7, 11	6, 10
6	5, 11	13, 14	12, 13	14, 10

题6的数据

参考文献

1. 薛毅, 陈立萍. R软件实用教程. 北京: 清华大学出版社, 2014.
2. 薛毅, 陈立萍. 统计建模与R软件(第二版). 北京: 清华大学出版社, 2021.
3. 张奕, 张彩旸. 应用统计学. 北京: 高等教育出版社, 2019.
4. Benjamini, Y. and Hochberg, Y. (1995), Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing, *Journal of the Royal Statistical Society: Series B (Methodological)*, Vol.57(1), 289–300.