

课程主页: <http://staff.ustc.edu.cn/~ynyang/2022>

第十五讲 方差分析

2022.12.9

线性性/可加性: 一个变量的效应
与其它变量的具体取值无关

因子变量/因素

因子/因素变量 (factor)

取值为类别的变量称为因子、因素、属性变量、类别变量，其各个取值称为水平 (level)，比如变量 Gender：取值为男、女，2个水平；Dose：取值高、中、低，3个水平（三个水平有次序，ordinal）。

单因素方差分析

根据一个因素（分类指标）将研究对象(随机)分为若干组。

例如：根据剂量分为高、中、低三组。因素：剂量，3个水平，如下表。

剂量	高	中	低

两因素方差分析

根据两个因素交叉分类 将研究对象随机分成双向表格(two-way)。

例如：两个药物A,B(两个因素)的用药剂量，称作因素A和因素B，各有三个水平(高、中、低)，交叉分类，共9组，如下表。

A \ B	高	中	低
高			
中			
低			

单因素方差分析 (one-way anova)

ANOVA: Analysis of variance

方差分析检验多个正态总体均值相等，是试验设计的重要分析方法。方差分析是线性回归模型中所有自变量为因子变量的特殊情形，虽然简单，但对于理解一般的回归分析思想也有帮助。

单因素方差分析模型：

随机化控制试验，研究对象被随机分配接受 K 种处理(*treatment*)中的一种。

假设第 k 组的响应服从（方差相同但均值可能不同的）正态分布：

$$y_{k1}, \dots, y_{kn_k} \text{ iid } \sim N(\mu_k, \sigma^2), \quad k = 1, \dots, K, \quad n = n_1 + \dots + n_K$$

零假设 $H_0: \mu_1 = \mu_2 = \dots = \mu_K$

这是多正态总体均值相等性检验，称为单因素方差分析模型。

因为随机化，在处理无效时所有响应的分布相同，在处理有效的情况下，只有均值可能有差异， N 和 σ^2 都相同。

处理	1	...	K
响应	$y_{11}, y_{12}, \dots, y_{1n_1}$...	$y_{K1}, y_{K2}, \dots, y_{Kn_K}$
平均	$\bar{y}_{1\cdot}$...	$\bar{y}_{K\cdot}$

$$\bar{y}_{k\cdot} = \sum_{i=1}^{n_k} y_{ki} / n_k,$$

$$\bar{y}_{\cdot\cdot} = \left(\sum_{k=1}^K \sum_{i=1}^{n_k} y_{ki} \right) / n,$$

$H_0: \mu_1 = \dots = \mu_K$ 的检验统计量应该基于 $\bar{y}_{1\cdot}, \dots, \bar{y}_{K\cdot}$ 之间的差异，如何度量 K 个数的差异？

我们可以两两比较并求平方和 $\sum_{i,j} (\bar{y}_{i\cdot} - \bar{y}_{j\cdot})^2$

也可以将每个 $\bar{y}_{i\cdot}$ 与总平均对比并求平方和 $\sum_i (\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot})^2$ ，事实上更合理的可能是加权平方和 $SS_B = \sum_i n_i (\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot})^2$ （因为第 i 组有 n_i 个样本， $\bar{y}_{i\cdot}$ 的方差与 $1/n_i$ 成正比），称为组间平方和。

线性模型形式

单因素方差分析模型（双下标形式）：
 $y_{ki} = \mu_k + \varepsilon_{ki}, \varepsilon_{ki} \text{ iid } \sim N(0, \sigma^2), i=1, \dots, n_k; k=1, \dots, K.$
 写成通常的线性模型形式：
 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{g}_1\mu_1 + \dots + \mathbf{g}_K\mu_K + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 I_n),$
 其中 $\mathbf{g}_1, \dots, \mathbf{g}_K$ 为各组的示性表示，相互正交， $X = (\mathbf{g}_1, \dots, \mathbf{g}_K)$. 具体如下：

$$\begin{array}{l}
 \text{第一组} \\
 \vdots \\
 y_{1n_1} \\
 \text{第二组} \\
 y_{21} \\
 \vdots \\
 y_{2n_2} \\
 \vdots \\
 y_{K1} \\
 \vdots \\
 y_{Kn_K} \\
 \text{第 } K \text{ 组}
 \end{array}
 \begin{pmatrix}
 y_{11} \\
 \vdots \\
 y_{1n_1} \\
 y_{21} \\
 \vdots \\
 y_{2n_2} \\
 \vdots \\
 y_{K1} \\
 \vdots \\
 y_{Kn_K}
 \end{pmatrix}
 =
 \begin{pmatrix}
 1 & 0 & \dots & 0 \\
 \vdots & \vdots & \dots & 0 \\
 1 & 0 & \dots & 0 \\
 0 & 1 & \dots & 0 \\
 \vdots & \vdots & & \\
 0 & 1 & \dots & 0 \\
 \vdots & \vdots & \ddots & \vdots \\
 0 & 0 & 0 & 1 \\
 \vdots & \vdots & \vdots & \vdots \\
 0 & 0 & 0 & 1
 \end{pmatrix}
 \begin{pmatrix}
 \mu_1 \\
 \mu_2 \\
 \vdots \\
 \mu_K
 \end{pmatrix}
 +
 \begin{pmatrix}
 \varepsilon_{11} \\
 \vdots \\
 \varepsilon_{1n_1} \\
 \varepsilon_{21} \\
 \vdots \\
 \varepsilon_{2n_2} \\
 \vdots \\
 \varepsilon_{K1} \\
 \vdots \\
 \varepsilon_{Kn_K}
 \end{pmatrix}
 =
 \begin{pmatrix}
 \mathbf{1}_{n_1} & \mathbf{0} & \dots & \mathbf{0} \\
 \mathbf{0} & \mathbf{1}_{n_2} & \dots & \mathbf{0} \\
 \vdots & \vdots & \dots & \vdots \\
 \mathbf{0} & \mathbf{0} & \dots & \mathbf{1}_{n_K}
 \end{pmatrix}
 \begin{pmatrix}
 \mu_1 \\
 \mu_2 \\
 \vdots \\
 \mu_K
 \end{pmatrix}
 + \boldsymbol{\varepsilon},$$

$(\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_K)$

注意 $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_K$ 是各个组的示性变量，为了在数学或计算机上处理双下标模型，我们将 K -水平因子变量转换为 K 个示性变量。

下面求 F 检验(检验 $q = K - 1$ 个自由参数等于0)

$$\hat{\mathbf{y}} = P_{\mathbf{g}_1} \mathbf{y} + \dots + P_{\mathbf{g}_K} \mathbf{y} = \mathbf{g}_1 \bar{y}_{1\bullet} + \dots + \mathbf{g}_K \bar{y}_{K\bullet} = (\bar{y}_{1\bullet}, \dots, \bar{y}_{1\bullet}, \bar{y}_{2\bullet}, \dots, \bar{y}_{2\bullet}, \dots, \bar{y}_{K\bullet}, \dots, \bar{y}_{K\bullet})^\top$$

H_0 下模型为: $\mathbf{y} = (\mathbf{g}_1 + \dots + \mathbf{g}_K) \mu_1 + \boldsymbol{\varepsilon} = \mathbf{1} \mu_1 + \boldsymbol{\varepsilon}$, $V_0 = L(\mathbf{1})$.

$$\hat{\mathbf{y}}_0 = P_{V_0} \mathbf{y} = (\bar{y}_{\bullet\bullet}, \dots, \bar{y}_{\bullet\bullet})^\top = \mathbf{1} \bar{y}_{\bullet\bullet}$$

$$\begin{aligned} SS_{\text{总}} \\ = SS_B + SS_W \end{aligned}$$

记号:

- 组间平方和: $SS_B = \|\hat{\mathbf{y}} - \hat{\mathbf{y}}_0\|^2 = \sum_{k=1}^K \sum_{i=1}^{n_k} (\bar{y}_{k\bullet} - \bar{y}_{\bullet\bullet})^2 = \sum_{k=1}^K n_k (\bar{y}_{k\bullet} - \bar{y}_{\bullet\bullet})^2$,
度量了各组的中心 $\bar{y}_{k\bullet}$, $k = 1, \dots, K$ 之间的差异 (回归平方和)。
- 组内平方和: $SS_W = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 = \sum_{k=1}^K \sum_{i=1}^{n_k} (y_{ki} - \bar{y}_{k\bullet})^2$,
度量了组内样本的差异 (残差平方和)。
- $SS_{\text{总}} = SS_B + SS_W$, B: Between - group, W: Within - group

$$\text{所以 } F = \frac{\|\hat{\mathbf{y}} - \hat{\mathbf{y}}_0\|^2 / q}{\|\mathbf{y} - \hat{\mathbf{y}}\|^2 / (n - p)} = \frac{SS_B / (K - 1)}{SS_W / (n - K)}$$

特例：两
样本t检验

两样本t检验是最简单的单因素方差分析方法。

$K = 2$ 时, 方差分析的 F 统计量: $F = (t)^2 \sim_{H_0} F_{1, n-2}$, 其中 t 为两样本 t 统计量

$$t = \frac{\bar{y}_{1\cdot} - \bar{y}_{2\cdot}}{\sqrt{(1/n_1 + 1/n_2)\hat{\sigma}^2}},$$

$$\text{其中 } \hat{\sigma}^2 = \left(\sum_{j=1}^{n_1} (y_{1j} - \bar{y}_{1\cdot})^2 + \sum_{j=1}^{n_2} (y_{2j} - \bar{y}_{2\cdot})^2 \right) / (n-2).$$

验证: $\bar{y}_{\cdot\cdot} = (n_1\bar{y}_{1\cdot} + n_2\bar{y}_{2\cdot}) / (n_1 + n_2)$, $n = n_1 + n_2$, 则

$$F = \frac{(n_1(\bar{y}_{1\cdot} - \bar{y}_{\cdot\cdot})^2 + n_2(\bar{y}_{2\cdot} - \bar{y}_{\cdot\cdot})^2) / (2-1)}{\left(\sum_{j=1}^{n_1} (y_{1j} - \bar{y}_{1\cdot})^2 + \sum_{j=1}^{n_2} (y_{2j} - \bar{y}_{2\cdot})^2 \right) / (n-2)} = \frac{n_1 n_2 (\bar{y}_{1\cdot} - \bar{y}_{2\cdot})^2}{n \hat{\sigma}^2} = \left(\frac{\bar{y}_{1\cdot} - \bar{y}_{2\cdot}}{\sqrt{(1/n_1 + 1/n_2)\hat{\sigma}^2}} \right)^2$$

附: 其它(单因素)多总体比较问题

非参数统计: 两总体Wilcoxon秩和检验, 多总体Kruskal-Wallis检验。

属性数据分析: $I \times J$ 列联表的Pearson卡方检验。

生存分析: log-rank检验。

例1 (npk数据). 随机在6块(block)地上种植豌豆, 每块地上4株。检验各块地上的产量是否相同. 我们需要检验6个地块的产量是否相同。 $H_0: \mu_1 = \dots = \mu_6$

block	1	2	3	4	5	6
产量	49.5,62.8, 46.8,57	59.8,58.5, 55.5,56	57.2,59, 53.2,56

```
> a = aov(yield ~ block, data = npk)
> summary( a )
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
block	5	343.3	68.66	2.318	0.0861 .
Residuals	18	533.1	29.61		

对比(contrast): 重新参数化

单因素方差分析模型中

$$y_{ki} \sim N(\mu_k, \sigma^2) \Leftrightarrow y_{ki} = \mu_k + \varepsilon_{ki}, \quad \varepsilon_{ki} \text{ iid} \sim N(0, \sigma^2) \quad (1)$$

零假设 $H_0: \mu_1 = \dots = \mu_K$ 是一个复合原假设。为了比较 K 个 μ , 可重新定义新的参数表达具体的感兴趣的问题, 称为对比(contrast)。

Treatment
Contrast

假设第一组(处理1)是安慰剂组(基准、baseline), 其它组使用各种不同药物(或同一药物的各种剂量), 那么我们关心的是各种药物或各种剂量相对于安慰剂的作用, 重新定义参数

$$\alpha_2 = \mu_2 - \mu_1, \dots, \alpha_K = \mu_K - \mu_1$$

称为各处理相对于基准的效应(effect)或处理对比(treatment contrast)。

$H_0: \mu_1 = \dots = \mu_K \Leftrightarrow H_0: \alpha_1 = \alpha_2 = \dots = \alpha_K = 0$, 这是一个简单原假设。

为了统一, 不妨记 $\alpha_1 \equiv 0$, 则 $\mu_k = \mu_1 + \alpha_k, k = 1, 2, \dots, K$. 模型(1)改写为

$$y_{ki} = \mu_1 + \alpha_k + \varepsilon_{ki}, \quad \varepsilon_{ki} \sim N(0, \sigma^2), \quad \text{约束 } \alpha_1 \equiv 0$$

Sum Contrast

如果没有一个公认的基准组，各个组对称，那么我们可以将每个处理与某个中心进行比较，比如取中心为 $\bar{\mu} = (\mu_1 + \dots + \mu_K)/K$ ，定义 K 个 sum对比：

$$\alpha_1 = \mu_1 - \bar{\mu}, \dots, \alpha_K = \mu_K - \bar{\mu},$$

注意有约束 $\alpha_1 + \dots + \alpha_K = 0$ 。原假设等价于 $H_0: \alpha_1 = \dots = \alpha_K = 0$ 。则单因素方差分析模型改写为

$$y_{ki} = \bar{\mu} + \alpha_k + \varepsilon_{ki}, \quad \varepsilon_{ki} \sim N(0, \sigma^2), \quad \text{约束 } \alpha_1 + \dots + \alpha_K = 0$$

一般的对比

若 $c_1 + \dots + c_K = 0$ ，称 $c_1\mu_1 + \dots + c_K\mu_K$ 是一个 对比(contrast)。

对比的定义方式取决于感兴趣的具体问题。比如 trt, sum contrast, 再如，如果 K 个组代表 K 个递增的剂量，假设我们关心相邻两种剂量的差异，则可定义

$$\alpha_k = \mu_{k+1} - \mu_k, \quad k = 1, 2, \dots, K - 1$$

无论何种约束，模型(1)都可改写为

$$y_{ki} = \mu + \alpha_k + \varepsilon_{ki}, \quad \varepsilon_{ki} \sim N(0, \sigma^2) \quad (2)$$

其中各组效应(对比) $\alpha_1, \dots, \alpha_K$ 满足一个线性约束。需要注意的是

- H_0 的F检验不受具体对比方式的影响，都如命题1所述。

引入对比概念的目的，一是具体问题的需要，二是方便对研究多因素方差分析。

- 双下标表示简明扼要，既表示了样本标号 i , 也表示了其所在的组号 k 。
- 但为了在数学上或计算机上处理，我们需要使用示性变量/哑变量表示分组（因子变量），如下页所示。

多下标的记录方式($a_{ij\dots k}$)称为数组array或张量tensor

对任何一个响应 y , 其所在组的标号记作 G , 方差分析模型写成线性模型:

$$y = \mu + \alpha_1 1_{(G=1)} + \dots + \alpha_K 1_{(G=K)} + \varepsilon \quad (3)$$

因为有个截距项, 诸 α 有一个线性约束, 所以实际上只需要 $K - 1$ 个哑变量, 比如trt contrast下 $\alpha_1 = 0$, 第一组的哑变量 $1_{(G=1)}$ 多余。

方差分析模型(3)的数据形式:

$$\begin{pmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ \vdots \\ y_{2n_2} \\ \vdots \\ y_{g1} \\ \vdots \\ y_{gn_g} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ \vdots & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & \dots & 1 \end{pmatrix} \begin{pmatrix} \mu_1 \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_K \end{pmatrix} + \varepsilon$$

← 如果是trt contrast约束, 即 $\alpha_1 = 0$, 实际计算时, 删除 α_1 及设计阵第二列。

因子的哑变量表示

哑变量表示 (dummy coding):

对于一个 K 水平因子变量, 取定其中一类为基准, 用其余 $K-1$ 个类的示性函数表示 K 个类。

例如, 因子变量 **group** 取值 **g1, g2, g3** (三个组), 以 **g1** 为基准, 定义

$$x_2 = 1_{(\text{group}=\text{g2})}, x_3 = 1_{(\text{group}=\text{g3})}$$

则二值变量 x_2 和 x_3 完全确定了 **group**:

$$\text{group} = \text{g1} \Leftrightarrow (x_2, x_3) = (0, 0);$$

$$\text{group} = \text{g2} \Leftrightarrow (x_2, x_3) = (1, 0);$$

$$\text{group} = \text{g3} \Leftrightarrow (x_2, x_3) = (0, 1).$$

group	x_2	x_3
g1	0	0
g2	1	0
g3	0	1

$\text{aov}(y \sim \text{group})$, 或 $\text{lm}(y \sim \text{group})$ 实际上拟合模型:

$$y = a + bx_2 + cx_3 + \varepsilon, \quad x_k = 1_{(\text{group}=k)}$$

其中 $\mu_1 = a, \mu_2 = a + b, \mu_3 = a + c \Leftrightarrow a = \mu_1, b = \mu_2 - \mu_1, c = \mu_3 - \mu_1$,
 b, c 分别称为第2,3组的效应, 它们是处理对比 (treatment contrast)。

配对t检验

配对设计

配对设计(paired design): 两个相像的研究对象（比如双胞胎），随机选取其中一人服用药物（处理），另一个服用安慰剂（对照）。

第*i*对样本的响应: (y_{1i}, y_{2i}) , 下标2表示处理,1表示对照。

$z_i = y_{2i} - y_{1i}$, 药效: $\delta = E(z_i)$, $H_0 : \delta = 0$

pair	1	2	...	<i>J</i>	平均
对照	y_{11}	y_{12}	...	y_{1J}	$\bar{y}_{1\cdot}$
处理	y_{21}	y_{22}	...	y_{2J}	$\bar{y}_{2\cdot}$
差	z_1	z_2	...	z_J	\bar{z}

同一个pair内的两个个体相似(匹配、相依), 将它们作为独立样本进行两样本t检验是错误的做法, 如何进行检验?

基于每对内的差 $z_j = y_{2j} - y_{1j}, j = 1, \dots, J$ 构造检验

成对t检验
(paired t-
test/depend
ent t-test)

配对样本 (y_{1i}, y_{2i}) , $i = 1, \dots, n$ iid, $z_i = y_{2i} - y_{1i}$, $\delta = E(z_i)$, $H_0 : \delta = 0$

假设 z_1, \dots, z_n iid $\sim N(\delta, \sigma^2)$, 成对t检验等于基于 z_i 的单正态总体均值检验:

$$t_{\text{pair}} = \frac{\bar{z}}{s_z/\sqrt{n}} = \frac{(\bar{y}_{2\cdot} - \bar{y}_{1\cdot})}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_{2i} - y_{1i} - (\bar{y}_{2\cdot} - \bar{y}_{1\cdot}))^2 / n}} \stackrel{H_0}{\sim} t_{n-1}$$

该检验称为成对 / 配对 t - 检验.

类比两样本t检验（比较处理组和对照组）：

$$t_{\text{two-sample}} = \frac{\bar{y}_{2\cdot} - \bar{y}_{1\cdot}}{\sqrt{(1/n + 1/n) \left(\frac{1}{2n-2} \sum_{j=1}^n (y_{1j} - \bar{y}_{1\cdot})^2 + \sum_{j=1}^n (y_{2j} - \bar{y}_{2\cdot})^2 \right)}}$$

t_{pair} 与 $t_{\text{two-sample}}$ 的分子相同，方差(分母)计算方式不同。

事实上, t_{pair} 可看作是如下方式得到的:

假设同一对的 y_{1i}, y_{2i} 共享一个共同的参数 μ_i (也可假设为随机变量) 用于刻画第 i 对两个研究对象的共性, 比如 μ_i 较大/小时, y_{1i}, y_{2i} 都较大/小。

$$\begin{cases} y_{1i} = \mu_i + \varepsilon_{1i} \\ y_{2i} = \mu_i + \delta + \varepsilon_{2i} \end{cases}, \quad \varepsilon_{1i}, \varepsilon_{2i}, i = 1, \dots, n \text{ iid} \sim N(0, \sigma^2)$$

令 $z_i = y_{2i} - y_{1i}$, 则

$$z_1, \dots, z_n \text{ iid} \sim N(\delta, 2\sigma^2),$$

消去了 μ_i , 因此将成对比较问题转化为单正态总体的均值检验 $H_0: \delta = 0$, 这正是上页的成对 t -检验的做法。

例2. McNemar检验：二值响应情形下的配对检验。

当响应变量 y_{ij} 为0,1变量（比如1表示治愈，否则为0），响应变量不可能服从正态分布，但我们依然可用成对 t 检验来检验处理的效应，只不过此时检验统计量在原假设下近似服从卡方分布。

<i>pair</i>	1	2	...	J
对照	y_{11}	y_{12}	...	y_{1J}
处理	y_{21}	y_{22}	...	y_{2J}

$$z_j = y_{2j} - y_{1j} = \begin{cases} 1 & y_{2j} = 1, y_{1j} = 0 \\ -1 & y_{2j} = 0, y_{1j} = 1 \\ 0 & y_{2j} = y_{1j} \end{cases}$$

原假设下 $\bar{z} \rightarrow 0$, $\sqrt{J}\bar{z} \rightarrow N(0, 2\sigma^2)$, 由Slusky引理, 知 $J \rightarrow \infty$ 时

$$t_{\text{pair}} = \sqrt{J}\bar{z} / \sqrt{(\sum_{j=1}^J z_j^2 - J\bar{z}^2) / (J-1)} \approx \sqrt{J}\bar{z} / \sqrt{\sum_{j=1}^J z_j^2 / J} = \sum_{j=1}^J z_j / \sqrt{\sum_{j=1}^J z_j^2} \xrightarrow{d} N(0, 1),$$

记 $b = \#\{j : y_{2j} = 1, y_{1j} = 0\}$, $c = \#\{j : y_{2j} = 0, y_{1j} = 1\}$ 分别为处理和对照响

应不同的pair的个数, 则 $\sum_{j=1}^J z_j = b - c$, $\sum_{j=1}^J z_j^2 = b + c$, $t_{\text{pair}} \approx (b - c) / \sqrt{b + c}$

配对二值(binary)数据可总结为下表, a 为两人响应都为0的pair的个数, b 为处理响应为1而对照响应为0的pair的个数, 等等。

		处理	
		0	1
对照	0	a	b
	1	c	d

$$t_{pair} \approx \frac{\sum_{j=1}^J z_j}{\sqrt{\sum_{j=1}^J z_j^2}} = \frac{(b-c)}{\sqrt{b+c}},$$

McNemar 检验: H_0 下, $X^2 = \frac{(b-c)^2}{b+c} \sim_{\text{近似}} \chi_1^2.$

- X^2 中没有 a, d : 两个研究对象响应都是0或者都是1的pair不提供处理效应的信息。
- McNemar检验是Cochran - Mantel - Haenszel(CMH)检验的特殊情况, CMH检验是综合多个列联表数据的独立性检验, 与两因素方差分析方法类似。

两因素方差分析

两因素方差分析中，决定分组的因素/因子有两个，可以都是处理（比如各种压力和各种温度组合下的工业试验），但通常一个因素是处理，另一个因素是区组，即区组设计。

随机区组设计 (Randomized block design)

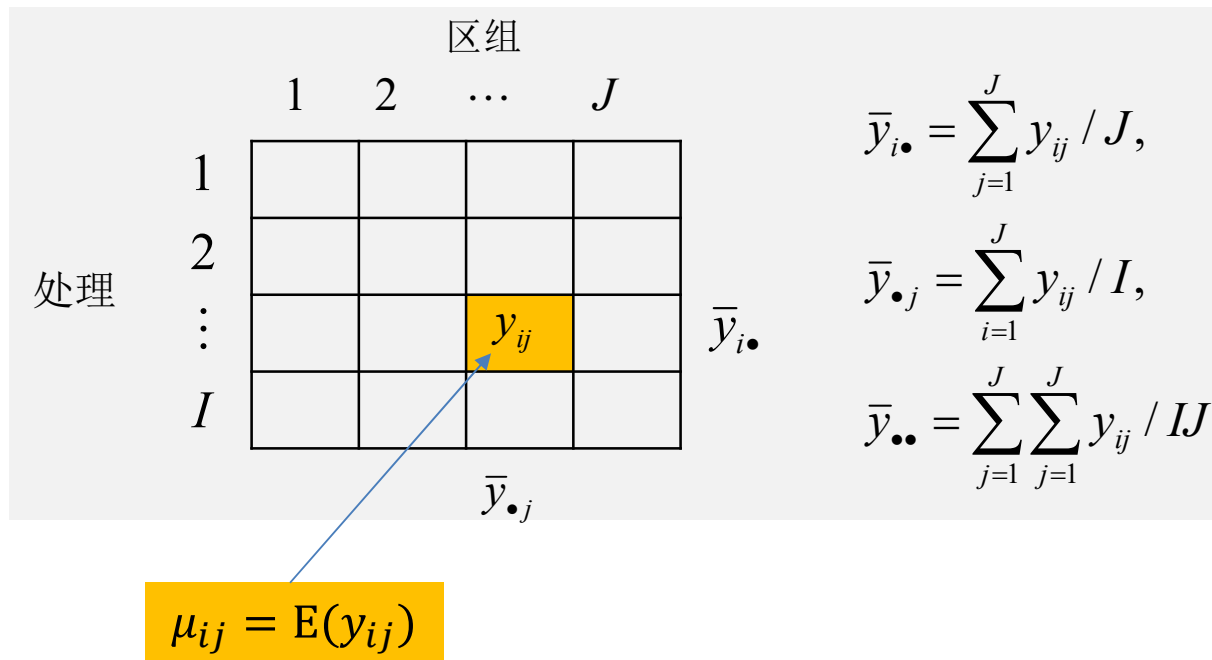
随机化试验设计中，当研究对象差异较大时，为了避免处理组与对照组出现严重的不均衡，通常先将研究对象分成内部个体相似的区组或层(blocks/strata), 在区组内随机化分配处理。

block: Fisher田间试验的地块;
strata: 生物、社会学常用

配对设计是最简单的区组设计，每一个pair就是一个区组。区组内的两个研究对象相似，但接受不同的处理。

假设两个因子各有 I, J 个水平，交叉分类得到 $I \times J$ 表，格子 (i, j) 内的观测数目为 n_{ij} 。

我们主要以每个格子内只有一次测量($n_{ij} = 1$)的情形为例（这也是最常见或最重要的情形），介绍两因素可加方差分析方法。假设 $y_{ij} \sim N(\mu_{ij}, \sigma^2)$ ，所有 y_{ij} 独立。



我们希望检验所有区组内，不同的处理的效果都相同：

$$H_0: \mu_{1j} = \mu_{2j} = \cdots = \mu_{Ij}, \quad 1 \leq j \leq J$$

		区组				
处理		1	...	j	...	J
1				μ_{1j}		
2				μ_{2j}		
⋮				⋮		
I				μ_{Ij}		

每个区组内只有 I 个样本，但有 $I+1$ 个参数，如何检验？

答案：假设 $\{\mu_{ij}\}$ 有简单结构/模型 - 假设各列之间存在共性(一致性)。

处理效应的一致性假设：

$\mu_{ij} - \mu_{i'j}$ 与所在区组 j 无关，只与处理水平 i, i' 有关

以 $I = J = 2$ 为例，我们说明效应一致性假设等价于可加性或线性性



所以处理效应一致性 \Leftrightarrow 区组效应一致性 \Leftrightarrow 可加性(两个因素互不干扰):

$$\mu_{21} = \mu_{11} + \alpha; \quad \mu_{12} = \mu_{11} + \beta \quad \mu_{22} = \mu_{21} + \beta = \mu_{11} + \alpha + \beta$$

即当一个因素固定而另一个因素从水平1到2时，均值增加 α 或 β 。当两个因素同时从水平1到2时，均值增加 $\alpha + \beta$ （可加）。

处理效应的一致性 \Leftrightarrow 区组效应的一致性 \Leftrightarrow 可加模型

一般的 $I \times J$ 情形类似。以处理1为基准(baseline), 假设处理的一致性: 处理 i ($1 \leq i \leq I$) 的效应 $\mu_{ij} - \mu_{1j}$ 与所在区组 j 无关, 记 $\alpha_i = \mu_{ij} - \mu_{1j}$, 称之为处理 i 的效应。

		区组1		区组j		
处理	1	μ_{11}		μ_{1j}		
	2	$\mu_{11} + \alpha_2$		$\mu_{1j} + \alpha_2$		
	\vdots	\vdots		\vdots		
	I	$\mu_{11} + \alpha_I$		$\mu_{1j} + \alpha_I$		

对任何一行 i (处理 i) , 区组1和区组j的均值之差,

$$\mu_{ij} - \mu_{i1} = (\mu_{1j} + \alpha_i) - (\mu_{11} + \alpha_i) = \mu_{1j} - \mu_{11} \triangleq \beta_j,$$
 与 i 无关 (即区组效应也有一致性), 记之为 β_j , 称为区组j的效应, 所以

$$\mu_{ij} = \mu_{i1} + \beta_j = \mu_{11} + \alpha_i + \beta_j$$

这称为可加效应模型 (处理和区组效应可加)。因此效应一致等价于可加模型。

两因素可加方差分析模型

假设处理 i ,区组 j 的响应 $y_{ij} \sim N(\mu_{ij}, \sigma^2)$, 其中

$$\mu_{ij} = \mu + \alpha_i + \beta_j \quad (\alpha's \text{ 和 } \beta's \text{ 各有一个约束})$$

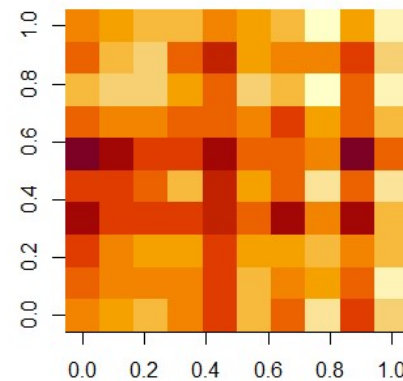
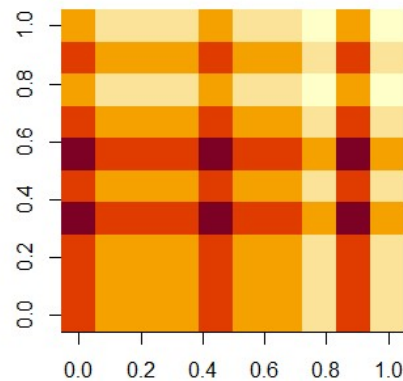
称为两因素可加方差分析模型 ($n_{ij} \equiv 1$ 情形)。

零假设为 $H_0: \alpha_1 = \alpha_2 = \dots = \alpha_I = 0$ 。

两因素可加模型是方差分析中最重要的模型和方法。

左图：可加模型 $\mu_{ij} = \mu + \alpha_i + \beta_j$ ，每两行之差为常数，每两列也是。

右图偏离可加，交互模型 $\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$ 。



可加方差分析模型 $y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}$ 是线性模型的特殊形式，
 例如 $I = J = 2$ 时，模型的矩阵-向量形式：

$$\begin{pmatrix} y_{11} \\ y_{21} \\ y_{12} \\ y_{22} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \boldsymbol{\varepsilon}$$

约束：

treatment contrast : $\alpha_1 = 0, \beta_1 = 0$

sum contrast : $\alpha_1 + \alpha_2 = 0, \beta_1 + \beta_2 = 0$

约束LS估计

可加模型下，参数个数为： $I+J-1$ ($I-1$ 个 α , $J-1$ 个 β , 一个 σ^2) $< n = IJ$ 。
参数可估。效应的LS估计依赖于具体的效应定义，我们仅演示sum contrast
约束情形（此时求解最简单），假设约束 $\Sigma\alpha_i = \Sigma\beta_j = 0$,

对误差平方和 $f = \sum_{i,j} (y_{ij} - \mu - \alpha_i - \beta_j)^2$, 求导得正则方程:

$$\begin{cases} \partial f / \partial \mu = \sum_i \sum_j (y_{ij} - \mu - \alpha_i - \beta_j) = 0, & IJ\mu = \sum_i \sum_j y_{ij} = IJ\bar{y}_{..} \Rightarrow \hat{\mu} = \bar{y}_{..} \\ \partial f / \partial \beta_j = \sum_i (y_{ij} - \mu - \alpha_i - \beta_j) = 0, & I\bar{y}_{.j} - I\mu - I\beta_j = 0 \Rightarrow \hat{\beta}_j = \bar{y}_{.j} - \bar{y}_{..} \\ \partial f / \partial \alpha_i = \sum_j (y_{ij} - \mu - \alpha_i - \beta_j) = 0, & J\bar{y}_{i.} - J\mu - J\alpha_i = 0 \Rightarrow \hat{\alpha}_i = \bar{y}_{i.} - \bar{y}_{..} \end{cases}$$

\Rightarrow 全模型下拟合值: $\hat{y}_{ij} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j = \bar{y}_{..} + (\bar{y}_{i.} - \bar{y}_{..}) + (\bar{y}_{.j} - \bar{y}_{..})$

零假设下，各处理(各行)无差异, $\tilde{\mu} = \bar{y}_{..}, \tilde{\beta}_j = \bar{y}_{.j} - \bar{y}_{..}$, 拟合值:

$$\tilde{y}_{ij} = \tilde{\mu} + \tilde{\beta}_j = \bar{y}_{..} + (\bar{y}_{.j} - \bar{y}_{..})$$

平方和分解

$$\begin{aligned} SS_{\text{总}} &= \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^I \sum_{j=1}^J \left((y_{ij} - \hat{y}_{ij}) + (\hat{y}_{ij} - \tilde{y}_{ij}) + (\tilde{y}_{ij} - \bar{y}_{..}) \right)^2 \\ &= \sum_{i=1}^I \sum_{j=1}^J \left((y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..}) + (\bar{y}_{i.} - \bar{y}_{..}) + (\bar{y}_{.j} - \bar{y}_{..}) \right)^2 \\ &= \sum_{i,j} (\bar{y}_{i.} - \bar{y}_{..})^2 + \sum_{i,j} (\bar{y}_{.j} - \bar{y}_{..})^2 + \sum_{i,j} (y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2 \stackrel{\Delta}{=} SS_{\text{处理}} + SS_{\text{区组}} + SS_{\text{组内}} \end{aligned}$$

- $SS_{\text{处理}} = \sum_{i=1}^I \sum_{j=1}^J (\bar{y}_{i.} - \bar{y}_{..})^2 = J \sum_{i=1}^I (\alpha_i - \bar{\alpha} + \bar{\varepsilon}_{i.} - \bar{\varepsilon}_{..})^2$, 处理的差异, 消除了区组效应
- $SS_{\text{组内}} = \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2$, 残差平方和, 消除了行、列效应

两因素方差分析F检验

命题2. I 个处理水平, J 个区组, 假设 (i, j) 水平下的测量 y_{ij} 满足可加模型

$$y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}, \quad \varepsilon_{ij} \text{ iid } \sim N(0, \sigma^2), \quad (\alpha's \text{ 和 } \beta's \text{ 各有一个约束})$$

$H_0: \alpha_1 = \dots = \alpha_I = 0$ 的F检验为:

$$F = \frac{\sum_{i=1}^I \sum_{j=1}^J (\bar{y}_{i.} - \bar{y}_{..})^2 / (J-1)}{\sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2 / (I-1)(J-1)} \stackrel{H_0}{\sim} F_{I-1, (I-1)(J-1)}$$

我们只需要验证自由度, $F = \frac{SS_{\text{处理}} / q}{SS_{\text{组内}} / (n-p)}$,

$q =$ 待检验的自由参数的个数 $= I - 1$ (注意 $\alpha's$ 有一个约束),

$p =$ 效应(回归系数)个数 $= (I-1) + (J-1) + 1 = I + J - 1$

$n =$ 样本量 $= IJ$, $n - p = IJ - I - J + 1 = (I-1)(J-1)$.

配对设计：
 $I = 2, n_{ij} \equiv 1$

当 $I = 2, n_{ij} = 1$ 时，F检验为： $F = J \bar{z}^2 / s_z^2 \sim F_{1, J-1}$ ，
其中 $t_{pair} = \sqrt{J} \bar{z} / s_z \sim t_{J-1}$ 为配对t检验统计量，
所以配对t检验是两因素方差分析的特殊情况。

验证：注意 $\bar{y}_{..} = (\bar{y}_{1.} + \bar{y}_{2.}) / 2 \Rightarrow \bar{y}_{1.} - \bar{y}_{..} = (\bar{y}_{1.} - \bar{y}_{2.}) / 2 = -\bar{y}_{2.} + \bar{y}_{..}$
 $SS_{处理} = J(\bar{y}_{1.} - \bar{y}_{..})^2 + J(\bar{y}_{2.} - \bar{y}_{..})^2 = J(\bar{y}_{1.} - \bar{y}_{2.})^2 / 2 = J \bar{z}^2 / 2$ ，
其中 $z_j = y_{2j} - y_{1j} \sim N(\alpha, 2\sigma^2)$ ， $\bar{z} = \bar{y}_{2.} - \bar{y}_{1.}$ ，其余略。

例3 ($I = J = 2$)

α : 处理效应, β : 区组效应

$$\begin{cases} y_{11} = \mu + \varepsilon_{11} \\ y_{21} = \mu + \alpha + \varepsilon_{21} \\ y_{12} = \mu + \beta + \varepsilon_{12} \\ y_{22} = \mu + \beta + \alpha + \varepsilon_{22} \end{cases}, \quad \varepsilon_{ij} \text{ iid } \sim N(0, \sigma^2)$$

		pair	
		1	2
对照	1	y_{11}	y_{12}
处理	2	y_{21}	y_{22}

令 $z_j = y_{2j} - y_{1j}, j = 1, 2$

$$H_0 : \alpha = 0, \quad F = \left(\frac{y_{21} - y_{11} + y_{22} - y_{12}}{y_{21} - y_{11} - y_{22} + y_{12}} \right)^2 = \left(\frac{z_1 + z_2}{z_1 - z_2} \right)^2 \sim_{H_0} F_{1,1}$$

$$t = \frac{z_1 + z_2}{z_1 - z_2} \sim_{H_0} t_1,$$

交互作用

仍以 2×2 情况为例，考察可加性或效应不一致情形下的模型。以处理1和区组1为基准。假设如果处理效应不一致，即 $\mu_{22} - \mu_{12} \neq \mu_{21} - \mu_{11}$ 。

	区组1	区组2
处理1	μ_{11}	μ_{12}
处理2	μ_{21}	μ_{22}

在基准区组1，假设处理效应为 $\alpha = \mu_{21} - \mu_{11}$



在基准处理组1，假设区组效应为 $\beta = \mu_{12} - \mu_{11}$

	区组1	区组2
处理1	μ_{11}	$\mu_{11} + \beta$
处理2	$\mu_{11} + \alpha$	μ_{22}

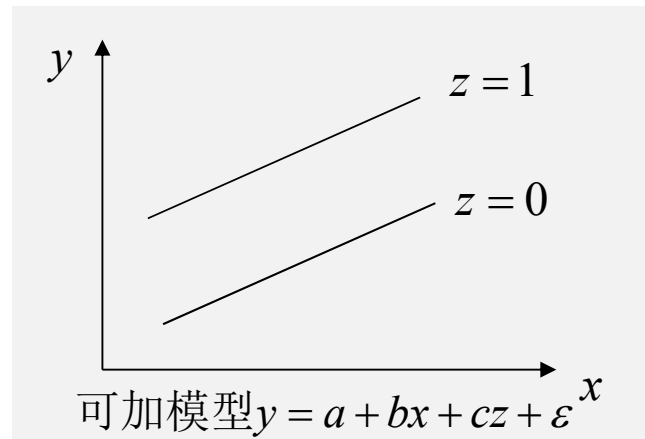
记效应之差 $\gamma = \gamma_{22} = (\mu_{22} - \mu_{12}) - (\mu_{21} - \mu_{11}) = \mu_{22} - \alpha - \beta - \mu_{11}$
则 $\mu_{22} = \mu_{11} + \alpha + \beta + \gamma$ 。

注意 $\gamma = 0$ 即 $\mu_{22} = \mu_{11} + \alpha + \beta$ 时是可加模型，所以 γ 代表了“当两个因素都从水平1变到2的时候， μ_{22} 偏离可加的程度”， γ 称为交互作用。

交互作用：一个因素的效应与另一个因素的具体水平有关

线性模型 \Leftrightarrow
效应一致性
 \Leftrightarrow 可加性

线性模型 $y = a + bx + cz + \varepsilon$ 意味着自变量效应的可加性：
不论自变量 z 取值如何， y 与 x 的关系(回归系数)都是常数 b .

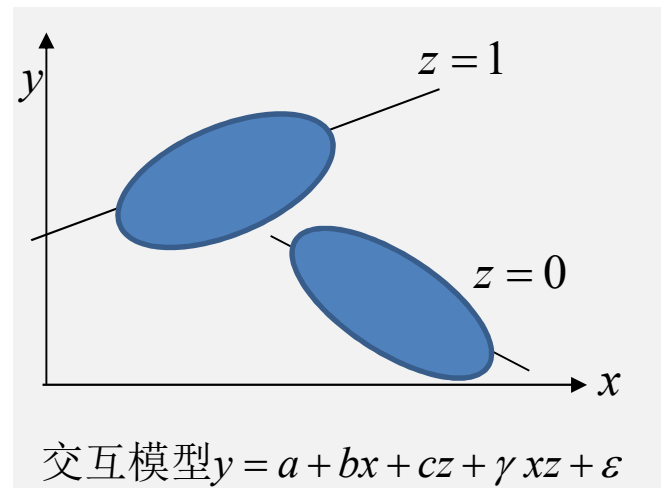


x 的效应与 z 无关：

$$z = 0 \text{ 时, } y = a + \underline{bx} + \varepsilon$$

$$z = 1 \text{ 时, } y = (a + c) + \underline{bx} + \varepsilon$$

交互作用：
偏离可加性



交互作用： x 的效应随 z 变化：

$$z = 0 \text{ 时, } y = a + \underline{bx} + \varepsilon$$

$$z = 1 \text{ 时, } y = (a + c) + \underline{(b + \gamma)x} + \varepsilon$$

附1: 平衡设计下的两因素可加方差分析 ($n_{ij} \equiv K$)

类似地, 对于一般的平衡设计 $n_{ij} \equiv K$, 有类似的结果
(对于 n_{ij} 不全相同的情形, 结果也类似, 在此不再赘述)

ANOVA F 检验

命题3: (两因素方差分析, 可加情形, 平衡设计: $n_{ij} \equiv K$)

原假设 $H_0: \alpha_1 = \alpha_2 = \dots = \alpha_I = 0$ 的 F 检验统计量为

$$F = \frac{SS_{\text{处理}} / (I-1)}{SS_{\text{组内}} / (n-I-J+1)} \stackrel{H_0}{\sim} F_{I-1, n-I-J+1}, \quad n = \sum_{i=1}^I \sum_{j=1}^J n_{ij} = IJK$$

其中 $SS_{\text{处理}} = \sum_{i,j,k} (\bar{y}_{i..} - \bar{y}_{...})^2$, $SS_{\text{组内}} = \sum_{i,j,k} (y_{ijk} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2$ 。

其中我们有方差分解公式:

$$\begin{aligned} SS_{\text{总}} &= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (y_{ijk} - \bar{y}_{...})^2 \stackrel{\Delta}{=} SS_{\text{处理}} + SS_{\text{区组}} + SS_{\text{组内}} \\ &= \sum_{i,j,k} (\bar{y}_{i..} - \bar{y}_{...})^2 + \sum_{i,j,k} (\bar{y}_{.j.} - \bar{y}_{...})^2 + \sum_{i,j,k} (y_{ijk} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2 \end{aligned}$$

附2. 交互作用模型

假设两个因素(x, z)都是0-1水平,组合(i, j)格子的均值 $\mu_{ij}, i, j = 0, 1$

饱和saturated模型/完全模型:

不对均值做任何结构假设

(4个参数: $\mu_{00}, \mu_{10}, \mu_{01}, \mu_{11}$).

$H_0: \mu_{00} = \mu_{10}, \mu_{01} = \mu_{11}$.

		z	
		0	1
x	0	μ_{00}	μ_{01}
	1	μ_{10}	μ_{11}

重新参数化(treatment contrast):

$$\mu = \mu_{00}, \quad \alpha_1 = \mu_{10} - \mu_{00}, \quad \beta_1 = \mu_{01} - \mu_{00}$$

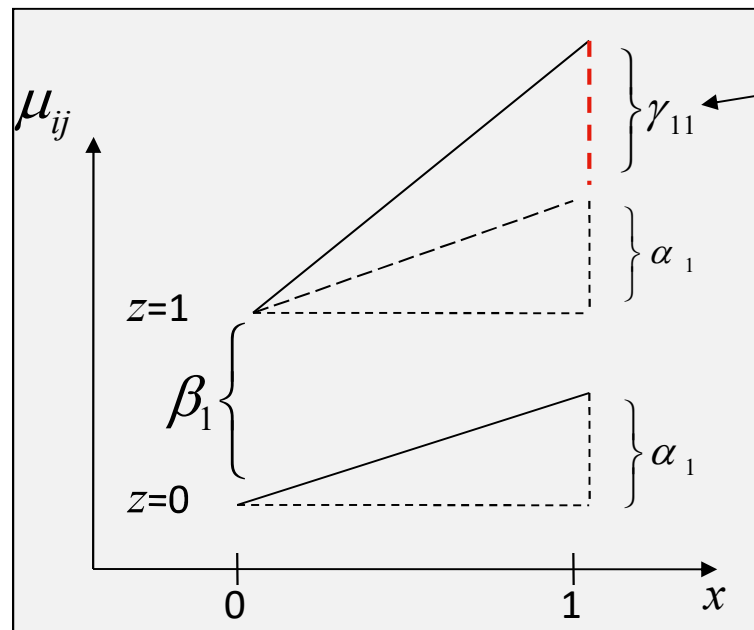
$$\begin{aligned} \gamma_{11} &= (\mu_{11} - \mu_{01}) - (\mu_{10} - \mu_{00}) \\ &= (\mu_{11} - \mu_{10}) - (\mu_{01} - \mu_{00}) \end{aligned}$$

		0	1
0	μ	$\mu + \beta_1$	
1	$\mu + \alpha_1$	$\mu + \alpha_1 + \beta_1 + \gamma_{11}$	

所以饱和模型可重写为交互作用模型：

$$\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$$

其中我们记 $\alpha_0 = 0, \beta_0 = 0, \gamma_{i0} = \gamma_{0j} = 0$ 。



交互效应

$$\begin{aligned} \gamma_{11} &= (\mu_{11} - \mu_{01}) - (\mu_{10} - \mu_{00}) \\ &= (\mu_{11} - \mu_{10}) - (\mu_{01} - \mu_{00}) \end{aligned}$$

$\gamma_{11} = 0 \Leftrightarrow$ 可加

两因素交互作用模型(一般情形)

设两个因子各 I, J 水平, 假设饱和模型/交互作用模型($n_{ij} \equiv K > 1$):

$$\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$$

其中参数 α, β, γ 个数分别是 $I-1, J-1, (I-1)(J-1)$ 个, 满足 $I+J+1$ 个约束。

treatment effect contrast 满足约束: $\alpha_1 = 0, \beta_1 = 0, \gamma_{1j} = \gamma_{i1} = 0,$

sum contrast 满足约束: $\sum_{i=1}^I \alpha_i = 0, \sum_{j=1}^J \beta_j = 0, \sum_{i=1}^I \gamma_{ij} = 0, \sum_{j=1}^J \gamma_{ij} = 0.$

sum contrast $\Leftrightarrow \{\mu_{ij}, i, j\}$ 重写参数化为

$$\mu = \sum \mu_{ij} / IJ, \alpha_i = \bar{\mu}_{i\cdot} - \mu = \sum \mu_{ij} / J - \mu,$$

$$\beta_j = \bar{\mu}_{\cdot j} - \mu = \sum \mu_{ij} / I - \mu, \gamma_{ij} = \mu_{ij} - \alpha_i - \beta_j + \mu.$$

LS估计

对误差平方和 $\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (y_{ijk} - \mu - \alpha_i - \beta_j - \gamma_{ij})^2$ 求导得正则方程:

$$\left\{ \begin{array}{l} IJK\mu + JK \sum_{i=1}^I \alpha_k + IK \sum_{j=1}^J \beta_j + K \sum_{i=1}^I \sum_{j=1}^J \gamma_{ij} = y_{\dots}, \\ JK\mu + JK\alpha_k + K \sum_{j=1}^J \beta_j + K \sum_{j=1}^J \gamma_{ij} = y_{i..}, \\ IK\mu + K \sum_{i=1}^I \alpha_i + IK\beta_j + K \sum_{i=1}^I \gamma_{ij} = \bar{y}_{.j.}, \\ K(\mu + \alpha_i + \beta_j + \gamma_{ij}) = y_{ij.}. \end{array} \right.$$

sum contrast 表示下, 参数满足约束 :

$$\sum_{i=1}^I \alpha_i = 0, \sum_{j=1}^J \beta_j = 0, \sum_{i=1}^I \gamma_{ij} = 0, \sum_{j=1}^J \gamma_{ij} = 0,$$

此时, LS估计具有简单形式 :

$$\begin{aligned} \hat{\mu} &= \bar{y}_{\dots}, & \hat{\alpha}_i &= \bar{y}_{i..} - \bar{y}_{\dots}, & \hat{\beta}_j &= \bar{y}_{.j.} - \bar{y}_{\dots}, \\ \hat{\gamma}_{ij} &= \bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{\dots} \end{aligned}$$

此时，平方和分解为：

$$\begin{aligned}
 SS_{\text{总}} &= \sum_{i,j,k} \left(\hat{\alpha}_i^2 + \hat{\beta}_j^2 + \hat{\gamma}_{ij}^2 + [y_{ijk} - (\hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + \hat{\gamma}_{ij})]^2 \right) \\
 &= \sum_{i,j,k} (\bar{y}_{i..} - \bar{y}_{...})^2 + \sum_{i,j,k} (\bar{y}_{.j.} - \bar{y}_{...})^2 + \sum_{i,j,k} (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2 \\
 &\quad + \sum_{i,j,k} (y_{ijk} - \bar{y}_{ij.})^2 \\
 &= SS_{\text{处理}} + SS_{\text{区组}} + SS_{\text{交互}} + SS_W
 \end{aligned}$$

可加性 检验

可加性检验：

$$H_0 : \gamma_{ij} = 0, \quad i = 1, \dots, I, j = 1, \dots, J \quad (\text{no interaction})$$

$$F_{\text{交互}} = \frac{SS_{\text{交互}} / (J-1)(I-1)}{SS_W / (IJK - IJ)} \sim_{H_0} F_{(I-1)(J-1), IJK - IJ}$$

$$\text{其中 } SS_{\text{交互}} = \sum_{i,j,k} (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2, \quad SS_W = \sum_{i,j,k} (y_{ijk} - \bar{y}_{ij.})^2$$

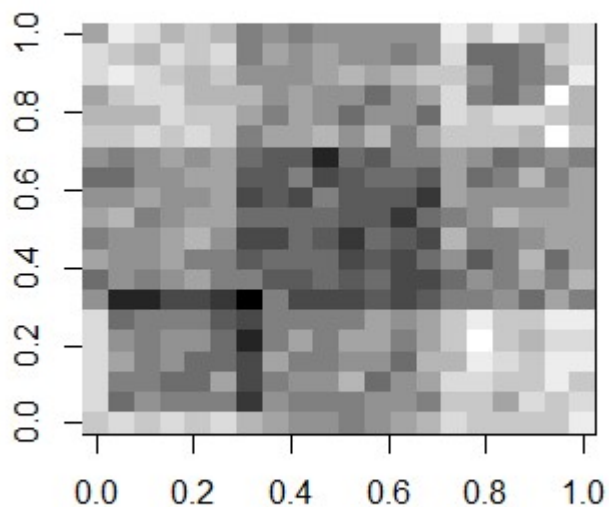
左图:原数据矩阵,

右图:消除可加模型(原数据减去行、列效应)后的残差,

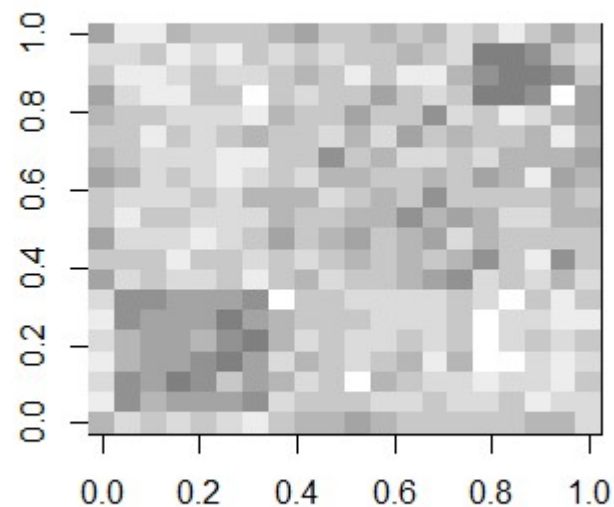
可加模型下,残差应该近似地服从 $(0, \sigma^2)$ 分布,

左下和右上角说明数据偏离可加模型,存在交互项。

$$y_{ij}$$



$$y_{ij} - \hat{\mu} - \hat{\alpha}_i - \hat{\beta}_j = y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{\cdot j} + \bar{y}_{\cdot\cdot}$$



Tukey 1-df 可加性检验

解决交互作用效应参数太多的问题

两因素, I, J 水平, 响应 $y_{ijk}, k = 1, 2, \dots, n_{ij}$ iid $\sim N(\mu_{ij}, \sigma^2)$, 其中
 $\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$.

$$H_0 : \gamma_{ij} = 0, i = 1, \dots, I, j = 1, \dots, J,$$

可加性检验的自由度为 $(I-1)(J-1), n - IJ$

- 当 I, J 较大时, 检验的参数个数/自由度为 $(I-1)(J-1)$ 太大, 功效不大。
- 另外一方面, 如果 $n_{ij} \equiv 1$ (没有重复测量), $RSS = 0$, 如何检验可加性?

Tukey 1-df 可加性检验解决了上述两个问题。

Tukey模型:

假设交互作用效应 $\gamma_{ij} = \lambda \alpha_i \beta_j$, 即 $\mu_{ij} = \mu + \alpha_i + \beta_j + \lambda \alpha_i \beta_j$

$$H_0 : \lambda = 0$$

算法 ($n_{ij} \equiv 1$) :

(1) 计算 $\hat{\mu} = \bar{y}_{..}$, $\hat{\alpha}_i = \bar{y}_{i.} - \bar{y}_{..}$, $\hat{\beta}_j = \bar{y}_{.j} - \bar{y}_{..}$, $\hat{\gamma}_{ij} = y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..}$

(2) 令 $w_{ij} = \hat{\alpha}_i \hat{\beta}_j$, 拟合 $y_{ij} = \hat{\mu} + \hat{\alpha}_i + \hat{\beta}_j + \lambda w_{ij} + \varepsilon_{ij}$, 即无截距的简单回归

$$\hat{\gamma}_{ij} = \lambda w_{ij} + \varepsilon_{ij}$$

$$\text{得LS估计: } \hat{\lambda} = \frac{\sum_i \sum_j w_{ij} y_{ij}}{\sum_i \sum_j w_{ij}^2} = \frac{\sum_i \sum_j \hat{\alpha}_i \hat{\beta}_j y_{ij}}{\sum_i \hat{\alpha}_i^2 \sum_j \hat{\beta}_j^2}$$

(3) 计算交互作用平方和, 残差平方和:

$$SS_{\text{int}} = \sum_i \sum_j (\hat{\lambda} \hat{\alpha}_i \hat{\beta}_j)^2, \quad SS_W = SS_{\text{总}} - \left(J \sum_i \hat{\alpha}_i^2 + I \sum_j \hat{\beta}_j^2 + SS_{\text{int}} \right).$$

(4) 计算 F 检验统计量: $F_{\text{Tukey}} = \frac{SS_{\text{int}}}{SS_W / (IJ - I - J)} \sim_{H_0} F_{1, IJ - I - J}$,

自由度为 **1**, $IJ - I - J$

附3：非参数检验和列联表检验

非参数检验和列联表检验与方差分析的对应关系

方法 总体分布 零分布	方差分析 正态分布 F分布	非参数检验 任何分布 卡方或置换	属性数据 多项分布 卡方或置换/Fisher
两总体（处理: 2水平）	两样本t检验	Wilcoxon秩和检验	Pearson卡方检验
多总体（处理: 多水平）	单因素方差分析	Kruskal-Wallis检验	Pearson卡方检验
多总体（处理×分层）	两因素方差分析	Friedman检验	CMH检验
配对 (2水平处理×分层)	配对t检验	Wilcoxon符号秩检验/Fisher符号检验	McNemar检验
无限总体	回归模型	非参数回归	Logistic回归/多项回归