# DOE 3

# Completely Randomized Design

## 1  Model and parameter estimation

### 1.1  Point estimation

Suppose that there are $r$ different treatments. We are interested in testing the equality of the $r$ population means of the treatments. Let $\mu_i$ be the population mean of $i$th treatment. Then the null and alternative hypotheses are

$$
\begin{aligned}
&H_0 : \mu_1 = \mu_2 = \cdots = \mu_r \\
&H_1 : \mu_i \neq \mu_j \quad \text{for at least one pair } (i,j)
\end{aligned}
\tag{1}
$$

where $i \neq j$. For the hypothesis testing above, suppose that we collected $r$ independent normal samples as below:

$$
\begin{aligned}
\text{Sample 1}: \ &Y_{11}, Y_{12}, \ldots, Y_{1n_1} \overset{iid}{\sim} N(\mu_1, \sigma^2) \\
\text{Sample 2}: \ &Y_{21}, Y_{22}, \ldots, Y_{2n_2} \overset{iid}{\sim} N(\mu_2, \sigma^2) \\
&\qquad \vdots \qquad\qquad\qquad \vdots \\
\text{Sample } r: \ &Y_{r1}, Y_{r2}, \ldots, Y_{rn_r} \overset{iid}{\sim} N(\mu_r, \sigma^2)
\end{aligned}
$$

Thus, it is reasonable to assume the following model

$$
Y_{ij} = \mu_i + \epsilon_{ij},
\tag{2}
$$

where $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, $i = 1, 2, \ldots, r$, and $j = 1, 2, \ldots, n_i$. This model is called the *one-way analysis of variance* (also known as the *one-way classification* or the *single-factor analysis*

*of variance*).

There is an alternative version of $Y_{ij} = \mu_i + \epsilon_{ij}$ in (2). Let $\mu_i = \mu + \tau_i$. That is, the $i$th treatment mean is equal to the *overall mean* (baseline) plus the $i$th treatment mean (effect) so that we have

$$Y_{ij} = \mu + \tau_i + \epsilon_{ij},$$

where we *assume* the following constraint

$$\sum_{i=1}^{r} n_i \tau_i = 0. \tag{3}$$

Note that the null hypothesis in (1) is equivalent to

$$\boxed{\begin{aligned} &H_0 : \tau_1 = \tau_2 = \cdots = \tau_r = 0 \\ &H_1 : \tau_i \neq 0 \quad \text{for at least one } i, \end{aligned}} \tag{4}$$

which is more simple than (1). A possible question can be *what if the constraint in (3) is not used?* In **Remark** 2, we will consider the other case.

One can estimate the parameters using the maximum likelihood estimation. The likelihood function is given by

$$L = \prod_{i=1}^{r} \prod_{j=1}^{n_i} \frac{1}{\sqrt{2\pi}\sigma} \exp\left( -\frac{\epsilon_{ij}^2}{2\sigma^2} \right) \tag{5}$$

and its log-likelihood is

$$\begin{aligned} \ell &= C - N \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^{r} \sum_{j=1}^{n_i} \epsilon_{ij}^2 \\ &= C - N \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2, \end{aligned} \tag{6}$$

where $N = \sum_{i=1}^{r} n_i$. Thus, maximizing the likelihood in (5) or the log-likelihood (6) is equivalent to minimizing $Q_2$ below

$$Q_2 = \sum_{i=1}^{r} \sum_{j=1}^{n_i} \epsilon_{ij}^2 = \sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2. \tag{7}$$

Differentiating $Q_2$ with respect to $\mu_i$ and setting it zero, we have

$$\frac{\partial Q_2}{\partial \mu_i} = 2 \sum_{j=1}^{n_i} (Y_{ij} - \mu_i) \cdot (-1) = (-2) \cdot \left( \sum_{j=1}^{n_i} Y_{ij} - n_i \mu_i \right) = 0.$$

Solving the above for $\mu_i$, we have

$$\hat{\mu}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} = \frac{1}{n_i} Y_{i\bullet} = \overline{Y}_{i\bullet},$$

where $Y_{i\bullet} = \sum_{j=1}^{n_i} Y_{ij}$.

We can also reparametrized $\mu_i$ with $\mu$ and $\tau_i$. Since $\mu_i = \mu + \tau_i$ and $\sum_{i=1}^{r} n_i \tau_i = 0$ due to (3), we have $\sum_{i=1}^{r} n_i \mu_i = \sum_{i=1}^{r} n_i \mu + \sum_{i=1}^{r} n_i \tau_i = N\mu$. Thus, we have $\mu = \sum_{i=1}^{r} n_i \mu_i / N$,

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{r} n_i \hat{\mu}_i = \frac{1}{N} \sum_{i=1}^{r} n_i \overline{Y}_{i\bullet} = \frac{1}{N} \sum_{i=1}^{r} \sum_{j=1}^{n_i} Y_{ij} = \overline{Y}_{\bullet\bullet} \tag{8}$$

and

$$\hat{\tau}_i = \hat{\mu}_i - \hat{\mu} = \overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet}, \tag{9}$$

where $\overline{Y}_{\bullet\bullet} = Y_{\bullet\bullet}/N$ and $Y_{\bullet\bullet} = \sum_{i=1}^{r} \sum_{j=1}^{n_i} Y_{ij}$.

It should be noted that the distribution of $\hat{\mu}$ is $\hat{\mu} \sim N(\mu, \sigma^2/N)$. For convenience, we denote $\hat{Y}_{ij} = \hat{\mu}_i = \hat{\mu} + \hat{\tau}_i$. Then we have

$$\hat{Y}_{ij} = \overline{Y}_{i\bullet}.$$

We also denote the residuals by $\hat{\epsilon}_{ij} = Y_{ij} - \hat{\mu}_i$. Then we have

$$\hat{\epsilon}_{ij} = Y_{ij} - \overline{Y}_{i\bullet}.$$

**Remark 1.** It should be noted that if the balanced design is used (that is, $n_1 = n_2 = \cdots = n_r = n$), then the constraint in (3) becomes

$$\sum_{i=1}^{r} \tau_i = 0. \tag{10}$$

$\triangle$

**Remark 2.** When we reparametrize $\mu_i$ with $\mu_i = \mu + \tau_i$, we used the constraint $\sum_{i=1}^{r} n_i \tau_i = 0$ in (3). Then a natural question arises: *What if the constraint $\sum_{i=1}^{r} \tau_i = 0$ is used?* If this constraint is used, we have $\sum_{i=1}^{r} \mu_i = \sum_{i=1}^{r} \mu + \sum_{i=1}^{r} \tau_i = r\mu$. Then we can estimate $\mu$ and $\tau_i$ as below.

$$\tilde{\mu} = \frac{1}{r} \sum_{i=1}^{r} \hat{\mu}_i = \frac{1}{r} \sum_{i=1}^{r} \overline{Y}_{i\bullet} = \overline{\overline{Y}}_{\bullet\bullet}$$

and

$$\tilde{\tau}_i = \hat{\mu}_i - \tilde{\mu} = \overline{Y}_{i\bullet} - \overline{\overline{Y}}_{\bullet\bullet}.$$

Notice that $\overline{\overline{Y}}_{\bullet\bullet}$ and $\overline{Y}_{\bullet\bullet}$ are not equal in general since

$$\overline{\overline{Y}}_{\bullet\bullet} = \frac{1}{r}\sum_{i=1}^{r}\frac{1}{n_i}\sum_{j=1}^{n_i}Y_{ij} \ \ \text{and} \ \ \overline{Y}_{\bullet\bullet} = \frac{1}{\sum_{i=1}^{r}n_i}\sum_{i=1}^{r}\sum_{j=1}^{n_i}Y_{ij}.$$

Thus, $\tilde{\mu}$ and $\tilde{\tau}_i$ are different from $\hat{\mu}$ in (8) and $\hat{\tau}_i$ (9), respectively. However, it is easily seen that $\overline{\overline{Y}}_{\bullet\bullet} = \overline{Y}_{\bullet\bullet}$ when the design is balanced.

Note that the distribution of $\overline{Y}_{i\bullet}$ is $\overline{Y}_{i\bullet} \overset{iid}{\sim} N(\mu+\tau_i, \sigma^2/n_i)$ due to (2). Then, under the constraint $\sum_{i=1}^{r}\tau_i = 0$, we have $\tilde{\mu} \sim N(\mu, \frac{\sigma^2}{r^2}\sum_{i=1}^{r}n_i^{-1})$ while $\hat{\mu} \sim N(\mu, \sigma^2/N)$. Thus, both $\hat{\mu}$ and $\tilde{\mu}$ are unbiased, but $\text{Var}(\hat{\mu}) \leq \text{Var}(\tilde{\mu})$ where the equality holds when the balanced design is used.          $\triangle$

## 1.2  Interval estimation

Recall that $Y_{ij} = \mu_i + \epsilon_{ij}$, where $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, $i = 1, 2, \ldots, r$, and $j = 1, 2, \ldots, n_i$. Then we have $Y_{i\bullet} \overset{iid}{\sim} N(n_i\mu_i, n_i\sigma^2)$ and thus $\overline{Y}_{i\bullet} \overset{iid}{\sim} N(\mu_i, \sigma^2/n_i)$. Standardizing $\overline{Y}_{i\bullet}$, we have

$$\frac{\overline{Y}_{i\bullet} - \mu_i}{\sqrt{\sigma^2/n_i}} \sim N(0, 1).$$

One can show that

$$\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij} - \overline{Y}_{i\bullet})^2 = \frac{1}{\sigma^2}\text{SSE} \sim \chi^2(N-r), \tag{11}$$

which will be detailed in Section 2.1.

---

**Definition 1.** The *covariance* of the random variables, $U$ and $V$ is defined by

$$\text{Cov}(U, V) = E[(U - \mu)(V - \nu)],$$

where $\mu = E(U)$ and $\nu = E(V)$.

---

**Lemma 1.**  *If $U_1, U_2, \ldots, U_n$ and $V_1, V_2, \ldots, V_m$ are random variables and $a_1, a_2, \ldots, a_n$ and $b_1, b_2, \ldots, b_m$ are constants, then we have*

$$\text{Cov}\left(\sum_{i=1}^{n}a_iU_i, \sum_{j=1}^{m}b_jV_j\right) = \sum_{i=1}^{n}\sum_{j=1}^{m}a_ib_j\text{Cov}(U_i, V_j).$$

---

*Proof.*  See Section 7.4 of Ross (2014).                                                                               $\square$

**Remark 3.**   Using Definition 1 and Lemma 1, it is easily seen that

$$\mathrm{Cov}(U_1, a_1) = 0$$

$$\mathrm{Cov}(U_1, U_1) = \mathrm{Var}(U_1)$$

$$\mathrm{Cov}(U_1, U_2) = \mathrm{Cov}(U_2, U_1)$$

$$\mathrm{Cov}(a_1 U_1, a_2 U_2) = a_1 a_2 \mathrm{Cov}(U_1, U_2)$$

$$\mathrm{Cov}(U_1 + a_1, U_2 + a_2) = \mathrm{Cov}(U_1, U_2)$$

$$\mathrm{Cov}(a_1 U_1 + a_2 U_2, b_1 V_1 + b_2 V_2) =$$

$$a_1 b_1 \mathrm{Cov}(U_1, V_1) + a_1 b_2 \mathrm{Cov}(U_1, V_2) + a_2 b_1 \mathrm{Cov}(U_2, V_1) + a_2 b_2 \mathrm{Cov}(U_2, V_2).$$

$$\triangle$$

---

**Lemma 2.**   *Under the assumption that $Y_{ij} = \mu_i + \epsilon_{ij}$ with $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, we have the following.*

*(a) $\mathrm{Cov}(\overline{Y}_{i\bullet}, Y_{ij} - \overline{Y}_{i\bullet}) = 0$.*

*(b) $\overline{Y}_{i\bullet}$ and $Y_{ij} - \overline{Y}_{i\bullet}$ are independent.*

*(c) $\overline{Y}_{i\bullet}$ and $\sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \overline{Y}_{i\bullet})^2$ are independent.*

---

*Proof.*   (a) Refer to Example 4e in Section 7.4 of Ross (2014).

(b) The proof of this part is based on (a) in this lemma. See Section 7.8.2 of Ross (2014). Also see Theorem 4.5-1 of Hogg et al. (2015).

(c) Note that this is very similar to Exercise 11.6 (b) of Casella and Berger (2002). For this proof, refer to Theorem 4.6.12 and Lemma 5.3.3 of Casella and Berger (2002) along with (b) in this lemma.   $\square$

We can easily show that $\mathrm{Cov}(\overline{Y}_{i\bullet}, Y_{ij} - \overline{Y}_{i\bullet}) = 0$ because

$$
\begin{aligned}
\mathrm{Cov}(\overline{Y}_{i\bullet}, Y_{ij} - \overline{Y}_{i\bullet}) &= \mathrm{Cov}(\overline{Y}_{i\bullet}, Y_{ij}) - \mathrm{Cov}(\overline{Y}_{i\bullet}, \overline{Y}_{i\bullet}) \\
&= \mathrm{Cov}\Big(\frac{1}{n_i} \sum_{j'=1}^{n_i} Y_{ij'}, Y_{ij}\Big) - \mathrm{Var}(\overline{Y}_{i\bullet}) \\
&= \frac{1}{n_i} \mathrm{Cov}(Y_{ij}, Y_{ij}) - \frac{\sigma^2}{n_i} = \frac{1}{n_i} \mathrm{Var}(Y_{ij}) - \frac{\sigma^2}{n_i} = 0.
\end{aligned}
$$

It should be noted that $\mathrm{Cov}(U, V) = 0$ does not guarantee the independence of $U$ and $V$ in general. However, as seen in Lemma 2, the condition $\mathrm{Cov}(U, V) = 0$ implies that independence of $U$ and $V$ especially for this one-way ANOVA model.

It is immediate from Lemma 2 (c) that $(\overline{Y}_{i\bullet} - \mu_i)/\sqrt{\sigma^2/n_i}$ and $\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij} - \overline{Y}_{i\bullet})^2/\sigma^2$ are also independent. Then we can Studentize $\overline{Y}_{i\bullet}$ as below

$$\frac{\frac{\overline{Y}_{i\bullet} - \mu_i}{\sqrt{\sigma^2/n_i}}}{\sqrt{\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij} - \overline{Y}_{i\bullet})^2/(N-r)}} \sim t(\mathrm{df} = N - r). \tag{12}$$

Note that we denote $\mathrm{SSE} = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij} - \overline{Y}_{i\bullet})^2$ and $\mathrm{MSE} = \mathrm{SSE}/(N-r)$ as will be given in (15) and (20), respectively. Using this MSE, we can rewrite (12) as

$$\frac{\overline{Y}_{i\bullet} - \mu_i}{\sqrt{\mathrm{MSE}/n_i}} \sim t(\mathrm{df} = N - r).$$

The endpoints for the interval estimation of $\mu_i$ with $100(1-\alpha)\%$ confidence level are given by

$$\boxed{\overline{Y}_{i\bullet} \ \pm \ t(1 - \tfrac{\alpha}{2}; N - r)\sqrt{\frac{\mathrm{MSE}}{n_i}},}$$

where $t(\gamma; \nu)$ is the *lower* $\gamma$th quantile of the $t$ distribution with $\nu$ degrees of freedom. This is also called the $100(1-\alpha)\%$ confidence interval of $\mu_i$.

We can also obtain the $100(1-\alpha)\%$ confidence interval of $\mu_\ell - \mu_m$ as follows. We have

$$\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet} \sim N\left(\mu_\ell - \mu_m, \sigma^2\left(\frac{1}{n_\ell} + \frac{1}{n_m}\right)\right)$$

and

$$\frac{(\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet}) - (\mu_\ell - \mu_m)}{\sqrt{\sigma^2\left(\frac{1}{n_\ell} + \frac{1}{n_m}\right)}} \sim N(0, 1).$$

Using (11), we have

$$\frac{(\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet}) - (\mu_\ell - \mu_m)}{\sqrt{\mathrm{MSE}\left(\frac{1}{n_\ell} + \frac{1}{n_m}\right)}} \sim t(\mathrm{df} = N - r). \tag{13}$$

Thus, the endpoints for the $100(1-\alpha)\%$ confidence interval of $\mu_\ell - \mu_m$ are given by

$$\boxed{\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet} \ \pm \ t(1 - \tfrac{\alpha}{2}; N - r)\sqrt{\mathrm{MSE}\left(\frac{1}{n_\ell} + \frac{1}{n_m}\right)}.} \tag{14}$$

# 2  Analysis of variance (ANOVA)

## 2.1  Decomposition of the total sum of squares

Since $Y_{ij} - \overline{Y}_{\bullet\bullet} = (Y_{ij} - \overline{Y}_{i\bullet}) + (\overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet}) = (\overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet}) + (Y_{ij} - \overline{Y}_{i\bullet})$, we have

$$\underbrace{\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2}_{\text{SS}_{\text{Total}}} = \underbrace{\sum_{i=1}^{r}\sum_{j=1}^{n_i}(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2}_{\text{SS}_{\text{Treatment}}} + \underbrace{\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2}_{\text{SS}_{\text{Error}}}$$

$$= \sum_{i=1}^{r} n_i(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2 + \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2 \qquad (15)$$

$$\text{SSTo} = \text{SStr} + \text{SSE}.$$

---

**Lemma 3.**  *Let $Z_i$ be* iid *standard normal random variables for $i = 1, 2, \ldots, \nu$ and $\sum_{i=1}^{\nu} Z_i^2 = V_1 + V_2 + \cdots + V_s$, where $V_j$ has $\nu_j$ degrees of freedom for $j = 1, 2, \ldots, s$ and $\nu_j > 0$. Then $V_j$ are independent with chi-squared random variables each with $\nu_j$ degrees of freedom if and only if $\nu = \nu_1 + \nu_2 + \cdots + \nu_s$.*

---

*Proof.*  See Cochran (1934) and Chapter 15 of Kendall and Stuart (1979).  $\square$

It is immediate from (2) that we have $Y_{ij} \overset{iid}{\sim} N(\mu_i, \sigma^2)$ for each $i$. Thus, we have $Y_{i\bullet} \overset{iid}{\sim} N(n_i\mu_i, n_i\sigma^2)$ and $Y_{\bullet\bullet} \sim N(\sum_{i=1}^{r} n_i\mu_i, N\sigma^2)$ which is equivalent to $Y_{\bullet\bullet} \sim N(N\mu, N\sigma^2)$ since $\sum_{i=1}^{r} n_i\mu_i = \sum_{i=1}^{r} n_i(\mu + \tau_i) = N\mu$ due to the constraint $\sum_{i=1}^{r} n_i\tau_i = 0$ from (3). Thus, $\overline{Y}_{i\bullet} \overset{iid}{\sim} N(\mu_i, \sigma^2/n_i)$ and $\overline{Y}_{\bullet\bullet} \sim N(\mu, \sigma^2/N)$. It should be emphasized that $\overline{Y}_{\bullet\bullet} \sim N(\mu, \sigma^2/N)$ due to the constraint $\sum_{i=1}^{r} n_i\tau_i = 0$, regardless of the null or alternative hypotheses, $H_0$ or $H_1$, in (1).

We have $(Y_{ij} - \mu)/\sigma \sim N(0, 1)$ under $H_0$, but $\sqrt{N}(\overline{Y}_{\bullet\bullet} - \mu)/\sigma \sim N(0, 1)$ regardless of $H_0$ or $H_1$. We have

$$\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\mu)^2 = \sum_{i=1}^{r}\sum_{j=1}^{n_i}\left\{(Y_{ij}-\overline{Y}_{\bullet\bullet}) + (\overline{Y}_{\bullet\bullet}-\mu)\right\}^2 = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2 + N(\overline{Y}_{\bullet\bullet}-\mu)^2$$

so that

$$\underbrace{\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\mu)^2}_{\chi^2(N)\text{ under }H_0} = \frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2 + \underbrace{\frac{1}{\sigma^2}N(\overline{Y}_{\bullet\bullet}-\mu)^2}_{\chi^2(1)}.$$

Using Lemma 3, we have

$$\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2 = \frac{1}{\sigma^2}\cdot\text{SSTo} \ \sim \ \chi^2(N-1) \text{ under } H_0. \tag{16}$$

We have $\sqrt{n_i}(\overline{Y}_{i\bullet}-\mu)/\sigma \sim N(0,1)$ under $H_0$, but $\sqrt{N}(\overline{Y}_{\bullet\bullet}-\mu)/\sigma \sim N(0,1)$ regardless of $H_0$ or $H_1$. We have

$$\sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\mu)^2 = \sum_{i=1}^{r}n_i\Big\{(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2+(\overline{Y}_{\bullet\bullet}-\mu)^2\Big\} = \sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2+N(\overline{Y}_{\bullet\bullet}-\mu)^2.$$

so that

$$\underbrace{\frac{1}{\sigma^2}\sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\mu)^2}_{\chi^2(r)\text{ under }H_0} = \frac{1}{\sigma^2}\sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2+\underbrace{\frac{1}{\sigma^2}N(\overline{Y}_{\bullet\bullet}-\mu)^2}_{\chi^2(1)}.$$

Then it is immediate from Lemma 3 that we have

$$\frac{1}{\sigma^2}\sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2 = \frac{1}{\sigma^2}\cdot\text{SStr} \ \sim \ \chi^2(r-1) \text{ under } H_0. \tag{17}$$

Again, we have SSTo = SStr + SSE from (15), that is,

$$\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2 = \sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2 + \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2,$$

which results in

$$\underbrace{\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2}_{\chi^2(N-1)\text{ under }H_0} = \underbrace{\frac{1}{\sigma^2}\sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2}_{\chi^2(r-1)\text{ under }H_0} + \frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2$$

due to (16) and (17). Thus, we have

$$\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2 \sim \chi^2(N-r) \text{ under } H_0. \tag{18}$$

We have shown that the statistic in (18) has a chi-squared distribution under $H_0$, but we can show that it has a chi-squared distribution regardless of $H_0$ and $H_1$ as below. We have

$$\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\mu_i)^2 = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2 + \sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\mu_i)^2.$$

Again, recall that we have $Y_{ij} \sim N(\mu_i,\sigma^2)$ and $\overline{Y}_{i\bullet} \sim N(\mu_i,\sigma^2/n_i)$ regardless of $H_0$ and $H_1$. Note that the distributions of the statistics in (16) and (17) were based on the $(Y_{ij}-\mu)/\sigma \sim N(0,1)$ under $H_0$, but $Y_{ij}$ is not restricted to $H_0$ here. We have

$$\underbrace{\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\mu_i)^2}_{\chi^2(N)} = \frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2 + \underbrace{\frac{1}{\sigma^2}\sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\mu_i)^2}_{\chi^2(r)}.$$

The above works under either $H_0$ and $H_1$ as aforementioned. Thus, we have

$$\frac{1}{\sigma^2}\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2 = \frac{1}{\sigma^2}\cdot\text{SSE} \sim \chi^2(N-r) \text{ under either } H_0 \text{ or } H_1. \tag{19}$$

In summary, we have

$$\boxed{\underbrace{\frac{1}{\sigma^2}\cdot\text{SSTo}}_{\chi^2(N-1)\text{ under }H_0} = \underbrace{\frac{1}{\sigma^2}\cdot\text{SStr}}_{\chi^2(r-1)\text{ under }H_0} + \underbrace{\frac{1}{\sigma^2}\cdot\text{SSE}}_{\chi^2(N-r)}.}$$

## 2.2 Expected mean square (EMS)

It is well known that the expected value of a chi-squared random variable is its degrees of freedom. Thus, using the above results, we have

$$E\Big(\frac{\text{SSTo}}{N-1}\Big) = E\Big(\frac{\text{SStr}}{r-1}\Big) = E\Big(\frac{\text{SSE}}{N-r}\Big) = \sigma^2 \text{ under } H_0.$$

Thus, $\text{SSTo}/(N-1)$, $\text{SStr}/(r-1)$ and $\text{SSE}/(N-r)$ are all unbiased for $\sigma^2$ under $H_0$. Then, what are the expected values of these quantities in general? First, for convenience, we denote

$$\text{MStr} = \frac{\text{SStr}}{r-1} \qquad \text{and} \qquad \text{MSE} = \frac{\text{SSE}}{N-r}, \tag{20}$$

which are called mean squares. The expected values of these mean squares are called EMS.

We can rewrite SSTo as

$$\text{SSTo} = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2 = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}^2 - 2\overline{Y}_{\bullet\bullet}Y_{ij} + \overline{Y}_{\bullet\bullet}^2) = \sum_{i=1}^{r}\sum_{j=1}^{n_i}Y_{ij}^2 - N\overline{Y}_{\bullet\bullet}^2.$$

Since $Y_{ij}\sim N(\mu_i,\sigma^2)$ and $\overline{Y}_{\bullet\bullet}\sim N(\mu,\sigma^2/N)$, we have $E(Y_{ij}^2) = \text{Var}(Y_{ij}) + \mu_i^2 = \sigma^2 + \mu_i^2$ and $E(\overline{Y}_{\bullet\bullet}^2) = \text{Var}(\overline{Y}_{\bullet\bullet}) + \mu^2 = \sigma^2/N + \mu^2$. Using these, we have

$$E(\text{SSTo}) = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(\sigma^2 + \mu_i^2) - (\sigma^2 + N\mu^2) = (N-1)\sigma^2 + \sum_{i=1}^{r}n_i\mu_i^2 - N\mu^2. \tag{21}$$

We also have

$$\sum_{i=1}^{r}n_i\mu_i^2 = \sum_{i=1}^{r}n_i(\mu+\tau_i)^2 = \sum_{i=1}^{r}n_i(\mu^2+\tau_i^2+2\mu\tau_i) = N\mu^2 + \sum_{i=1}^{r}n_i\tau_i^2 + 2\mu\sum_{i=1}^{r}n_i\tau_i,$$

where the term, $\sum_{i=1}^{r}n_i\tau_i$, is zero by (3) so that we have

$$\sum_{i=1}^{r}n_i\mu_i^2 = N\mu^2 + \sum_{i=1}^{r}n_i\tau_i^2. \tag{22}$$

Substituting (22) into (21), we have

$$E(\text{SSTo}) = (N-1)\sigma^2 + \sum_{i=1}^{r} n_i \tau_i^2 \;\; \text{and} \;\; E\Big(\frac{\text{SSTo}}{N-1}\Big) = \sigma^2 + \frac{\sum_{i=1}^{r} n_i \tau_i^2}{N-1}.$$

As expected, $E\big(\text{SSTo}/(N-1)\big)$ becomes $\sigma^2$ under $H_0$.

We rewrite SStr as

$$\text{SStr} = \sum_{i=1}^{r} n_i (\overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet})^2 = \sum_{i=1}^{r} n_i \overline{Y}_{i\bullet}^2 - N \overline{Y}_{\bullet\bullet}^2.$$

Since $\overline{Y}_{i\bullet} \sim N(\mu_i, \sigma^2/n_i)$, we have $E(\overline{Y}_{i\bullet}^2) = \text{Var}(\overline{Y}_{i\bullet}) + \mu_i^2 = \sigma^2/n_i + \mu_i^2$. Using this along with (22) and $E(\overline{Y}_{\bullet\bullet}^2) = \sigma^2/N + \mu^2$, we have

$$E(\text{SStr}) = \sum_{i=1}^{r}(\sigma^2 + n_i \mu_i^2) - (\sigma^2 + N\mu^2) = (r-1)\sigma^2 + \sum_{i=1}^{r} n_i \mu_i^2 - N\mu^2 = (r-1)\sigma^2 + \sum_{i=1}^{r} n_i \tau_i^2,$$

which results in

$$\boxed{E(\text{MStr}) = \sigma^2 + \frac{\sum_{i=1}^{r} n_i \tau_i^2}{r-1}.}$$

This quantity also becomes $\sigma^2$ under $H_0$.

Next, we rewrite SSE as

$$\text{SSE} = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij} - \overline{Y}_{i\bullet})^2 = \sum_{i=1}^{r}\sum_{j=1}^{n_i} Y_{ij}^2 - \sum_{i=1}^{r} n_i \overline{Y}_{i\bullet}^2$$

Using $Y_{ij} \sim N(\mu_i, \sigma^2)$ and $\overline{Y}_{i\bullet} \sim N(\mu_i, \sigma^2/n_i)$, we have

$$E(\text{SSE}) = \sum_{i=1}^{r}\sum_{j=1}^{n_i}(\sigma^2 + \mu_i^2) - \sum_{i=1}^{r} n_i(\sigma^2/n_i + \mu_i^2) = (N-r)\sigma^2,$$

which results in

$$\boxed{E(\text{MSE}) = \sigma^2.}$$

Note that the above works under $H_0$ and $H_1$ as well.

### ANOVA Decomposition

| Source | SS | df | MS | $F$ | EMS |
|---|---|---|---|---|---|
| Treatment | SStr | $r-1$ | $\text{MStr} = \text{SStr}/(r-1)$ | $F = \dfrac{\text{MStr}}{\text{MSE}}$ | $\sigma^2 + \dfrac{\sum_{i=1}^{r} n_i \tau_i^2}{r-1}$ |
| Error | SSE | $N-r$ | $\text{MSE} = \text{SSE}/(N-r)$ | | $\sigma^2$ |
| Total | SSTo | $N-1$ | | | |

In summary, we recall the decomposition of the total sum of squares in (15):

$$\underbrace{\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{\bullet\bullet})^2}_{\text{SSTo}} = \underbrace{\sum_{i=1}^{r}n_i(\overline{Y}_{i\bullet}-\overline{Y}_{\bullet\bullet})^2}_{\text{SStr}} + \underbrace{\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij}-\overline{Y}_{i\bullet})^2}_{\text{SSE}} \quad (23)$$

# 3  One-way ANOVA as regression

## 3.1  Cell-means coding

A very important application of regression analysis involves a list of predictors which can include categorical variables for treatments. The categorical variables are also called *indicator* or *dummy* variables or qualitative variables. We re-analyze the above one-way ANOVA model using the linear regression model. For notational brevity, we omit $ij$ index and we then have

$$Y = \beta_1 Z_1 + \beta_2 Z_2 + \cdots + \beta_r Z_r + \epsilon \quad (24)$$

where $\epsilon \overset{iid}{\sim} N(0,\sigma^2)$. We define $Z_i$ as follows:

$$Z_i = \begin{cases} 1 & : \quad \text{for the } i\text{th treatment} \\ 0 & : \quad \text{otherwise} \end{cases},$$

where $i = 1, 2, \ldots, r$. That is, the dummy variables $Z_i$ are coded as below. This scheme is called *cell-means coding*

<div align="center">

Cell-means coding

| Treatment | $Z_1$ | $Z_2$ | $Z_3$ | $\cdots$ | $Z_r$ |
|:---:|:---:|:---:|:---:|:---:|:---:|
| $i=1$ | 1 | 0 | 0 | $\cdots$ | 0 |
| $i=2$ | 0 | 1 | 0 | $\cdots$ | 0 |
| $i=3$ | 0 | 0 | 1 | $\cdots$ | 0 |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $i=r$ | 0 | 0 | 0 | $\cdots$ | 1 |

</div>

It is easy to see that the model in (24) implies that

$$\mu_1 = E(Y_{1j}) = \beta_1$$

$$\mu_2 = E(Y_{2j}) = \beta_2$$

$$\vdots = \quad \vdots \quad = \vdots$$

$$\mu_r = E(Y_{rj}) = \beta_r$$

where $Y_{ij} = \mu_i + \epsilon_{ij}$ and $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$ again. This clearly shows that the linear regression can handle the one-way ANOVA model using the cell-means coding.

Suppose that we stack the responses $(Y_i)$ from the $r$ samples in such a way that the first sample comes first, the second sample comes next, and the last $r$th sample comes last. Then the response vector of length $N = \sum_{i=1}^{r} n_i$, the data matrix (or design matrix) and the vector of the error terms are given by

$$
\underset{N\times 1}{\mathbf{Y}} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_{n_1} \\ Y_{n_1+1} \\ \vdots \\ Y_{n_1+n_2} \\ \vdots \\ \vdots \\ Y_{N-n_r+1} \\ \vdots \\ Y_N \end{bmatrix} = \begin{bmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ \vdots \\ Y_{r1} \\ \vdots \\ Y_{rn_r} \end{bmatrix}, \underset{N\times r}{\mathbf{Z}} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}, \underset{N\times 1}{\boldsymbol{\epsilon}} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_{n_1} \\ \epsilon_{n_1+1} \\ \vdots \\ \epsilon_{n_1+n_2} \\ \vdots \\ \vdots \\ \epsilon_{N-n_r+1} \\ \vdots \\ \epsilon_N \end{bmatrix} = \begin{bmatrix} \epsilon_{11} \\ \vdots \\ \epsilon_{1n_1} \\ \epsilon_{21} \\ \vdots \\ \epsilon_{2n_2} \\ \vdots \\ \vdots \\ \epsilon_{r1} \\ \vdots \\ \epsilon_{rn_r} \end{bmatrix}.
$$

Using the response vector and the data matrix, we can rewrite the linear regression model as

$$\boxed{\mathbf{Y} = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\epsilon}.}$$

To estimate the regression parameters, we use the least squares method which is equivalent to the maximum likelihood method under the normal distribution assumption. It is easily seen that the sum of squared errors is given by

$$Q_2 = \sum_{i=1}^{r} \sum_{j=1}^{n_i} \epsilon_{ij}^2 = \sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \beta_i)^2,$$

which is essentially the same as the (7). Differentiating $Q_2$ with respect to $\beta_i$ and solving the estimating equations, we have

$$\hat{\beta}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} = \frac{1}{n_i} Y_{i\bullet} = \overline{Y}_{i\bullet},$$

where $Y_{i\bullet} = \sum_{j=1}^{n_i} Y_{ij}$ again.

**Remark 4.** It is easy to show that

$$\mathbf{Z}'\mathbf{Z} = \text{diag}(n_1, n_2, \ldots, n_r) = \begin{bmatrix} n_1 & 0 & \ldots & 0 \\ 0 & n_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \ldots & n_r \end{bmatrix}.$$

Thus, we have $(\mathbf{Z}'\mathbf{Z})^{-1} = \text{diag}(1/n_1, 1/n_2, \ldots, 1/n_r)$. Let $\mathbf{1}_n$ be a $n$-dimensional vector with all the elements being ones. Then $\mathbf{J}_n = \mathbf{1}_n \mathbf{1}_n'$ is a $n \times n$ square matrix of ones, which is often called *all-ones matrix*. Then the hat matrix for the cell-means coding is given by

$$\mathbf{H} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}' = \text{diag}\left(\frac{1}{n_1}\mathbf{J}_{n_1}, \frac{1}{n_2}\mathbf{J}_{n_2}, \ldots, \frac{1}{n_r}\mathbf{J}_{n_r}\right).$$

$\triangle$

**Remark 5.** The cell-means coding has nice features. However, it is rarely used in practice because of the following reasons.

1. The data matrix $\mathbf{Z}$ does not include $\mathbf{1} = (1, 1, \ldots, 1)'$ in the column which implies that the intercept (grand mean) is not used. Thus, a usual way of decomposing the total sum of squares (SSTo) does not work.

2. Based on the regression model in (24), one can test the hypothesis that a cell mean is zero in such a way that

    $$H_0 : \beta_i = 0,$$

    whose usual testing statistic is given by

    $$T = \frac{\hat{\beta}_i - 0}{\sqrt{\text{MSE} \cdot (\mathbf{Z}'\mathbf{Z})_{ii}^{-1}}} = \frac{\hat{\beta}_i - 0}{\sqrt{\text{MSE}/n_i}}.$$

    Unlike the conventional regression applications, it does not often make sense to perform the hypothesis above in real-world applications of the design of experiments. Usually, it is more interesting to test whether $\tau_i = 0$ instead of $\beta_i = 0$.

For these reasons, the cell-means coding is not widely-used in practice. However, the use of the cell-means coding along with the regression method present an idea of how to incorporate the regression method into the design of experiments. $\triangle$

## 3.2  Effect coding

The effect is another popular coding scheme to incorporate the regression method into the design of experiments. We re-analyze the above one-way ANOVA model using the linear regression model *with the intercept.* For notational brevity, we omit $ij$ index and we then have

$$Y = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2 + \cdots + \beta_{r-1} Z_{r-1} + \epsilon \tag{25}$$

where $\epsilon \overset{iid}{\sim} N(0, \sigma^2)$. For the $i$th treatment where $i = 1, 2, \ldots, r-1$, we set up $Z_i = 1$ and all other dummy variables are zero, but we set $Z_i = -n_i/n_r$ for the $r$th treatment. Then we have

$$Z_i = \begin{cases} 1 & : & \text{for the } i\text{th treatment} \\ -n_i/n_r & : & \text{for the } r\text{th treatment} \\ 0 & : & \text{otherwise} \end{cases},$$

where $i = 1, 2, \ldots, r-1$. This scheme is called *effect coding.* Note that, including the intercept $\beta_0$, there are $r$ regression coefficients all told. For convenience, let $Z_0 = 1$ which is considered as the predictor with $\beta_0$. That is, the dummy variables $Z_i$ are coded as below.

<div align="center">

Effect coding

| Treatment | $Z_0$ | $Z_1$ | $Z_2$ | $\cdots$ | $Z_{r-1}$ |
|:---:|:---:|:---:|:---:|:---:|:---:|
| $i = 1$ | 1 | 1 | 0 | $\cdots$ | 0 |
| $i = 2$ | 1 | 0 | 1 | $\cdots$ | 0 |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $i = r-1$ | 1 | 0 | 0 | $\cdots$ | 1 |
| $i = r$ | 1 | $-n_1/n_r$ | $-n_2/n_r$ | $\cdots$ | $-n_{r-1}/n_r$ |

</div>

It is easy to see that the model in (25) implies that

$$\begin{aligned}
\mu_1 &= E(Y_{1j}) & &= \beta_0 + \beta_1 \\
\mu_2 &= E(Y_{2j}) & &= \beta_0 + \beta_2 \\
\vdots &= \vdots & &= \vdots \\
\mu_{r-1} &= E(Y_{r-1,j}) & &= \beta_0 + \beta_{r-1} \\
\mu_r &= E(Y_{rj}) & &= \beta_0 - \frac{n_1}{n_r}\beta_1 - \frac{n_2}{n_r}\beta_2 - \cdots - \frac{n_{r-1}}{n_r}\beta_{r-1}
\end{aligned}$$

where $Y_{ij} = \mu_i + \epsilon_{ij}$ and $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$ again.

Again, suppose that we stack the responses $(Y_i)$ from the $r$ samples in such a way that the first sample comes first, the second sample comes next, and the last $r$th sample comes last. Then the response vector of length $N = \sum_{i=1}^r n_i$ and the data or design matrix are given by

$$
\underset{N \times 1}{\mathbf{Y}} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_{n_1} \\ Y_{n_1+1} \\ \vdots \\ Y_{n_1+n_2} \\ \vdots \\ \vdots \\ Y_{N-n_r+1} \\ \vdots \\ Y_N \end{bmatrix} = \begin{bmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ \vdots \\ Y_{r1} \\ \vdots \\ Y_{rn_r} \end{bmatrix}, \underset{N \times r}{\mathbf{Z}} = \begin{bmatrix} 1 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots & \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & & \vdots \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ 1 & -n_1/n_r & -n_2/n_r & \cdots & -n_{r-1}/n_r \\ \vdots & \vdots & & & \vdots \\ 1 & -n_1/n_r & -n_2/n_r & \cdots & -n_{r-1}/n_r \end{bmatrix}.
$$

Using the response vector and the data matrix, we can rewrite the linear regression model as

$$\boxed{\mathbf{Y} = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\epsilon}.}$$

Analogous with the cell-means coding case, we can estimate the regression parameters using the least squares method. Then the sum of squared errors is given by

$$
\begin{aligned}
Q_2 &= \sum_{i=1}^r \sum_{j=1}^{n_i} \epsilon_{ij}^2 \\
&= \sum_{i=1}^{r-1} \sum_{j=1}^{n_i} \epsilon_{ij}^2 + \sum_{j=1}^{n_i} \epsilon_{rj}^2 \\
&= \sum_{i=1}^{r-1} \sum_{j=1}^{n_i} (Y_{ij} - \beta_0 - \beta_i)^2 + \sum_{j=1}^{n_i} \left( Y_{rj} - \beta_0 + \frac{n_1}{n_r}\beta_1 + \frac{n_2}{n_r}\beta_2 + \cdots + \frac{n_{r-1}}{n_r}\beta_{r-1} \right)^2.
\end{aligned}
$$

Let $\beta_r = -(n_1\beta_1 + n_2\beta_2 + \cdots + n_{r-1}\beta_{r-1})/n_r$. Then the minimization of the above $Q_2$ is equivalent to minimizing it

$$Q_2^* = \sum_{i=1}^r \sum_{j=1}^{n_i} (Y_{ij} - \beta_0 - \beta_i)^2 \quad \text{subject to} \quad \sum_{i=1}^r n_i\beta_i = 0.$$

The auxiliary function with Lagrange multiplier $\lambda$ is

$$\Psi = \sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \beta_0 - \beta_i)^2 - \lambda \sum_{i=1}^{r} n_i \beta_i.$$

Differentiating $\Psi$ with respect to $\beta_0$ and setting it to zero, we have

$$\sum_{i=1}^{r} \sum_{j=1}^{n_i} (Y_{ij} - \beta_0 - \beta_i) = 0,$$

which results in

$$Y_{\bullet\bullet} - N\beta_0 - \sum_{i=1}^{r} n_i \beta_i = Y_{\bullet\bullet} - N\beta_0 = 0.$$

Thus, we have

$$\hat{\beta}_0 = \overline{Y}_{\bullet\bullet}. \tag{26}$$

Next, differentiating $\Psi$ with respect to $\beta_i$ and setting it to zero, we have

$$\sum_{j=1}^{n_i} (Y_{ij} - \beta_0 - \beta_i) = 0,$$

which results in

$$Y_{i\bullet} - n_i \beta_0 - n_i \beta_i = 0.$$

Solving the above for $\beta_i$ and substituting $\hat{\beta}_0 = \overline{Y}_{\bullet\bullet}$ into the above, we have

$$\hat{\beta}_i = \overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet}. \tag{27}$$

Comparing (26) and (27) with (8) and (9), respectively, it is easily seen that

$$\hat{\beta}_0 = \hat{\mu} \quad \text{and} \quad \hat{\beta}_i = \hat{\tau}_i.$$

**Remark 6.** Unlike the the cell-means coding case, it is quite complex to obtain $\mathbf{Z}'\mathbf{Z}$ which is certainly different from the cell-means coding case. However, the hat matrix for the effect coding is the same as that for the cell-means coding which is again given by

$$\mathbf{H} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}' = \text{diag}\left(\frac{1}{n_1}\mathbf{J}_{n_1}, \frac{1}{n_2}\mathbf{J}_{n_2}, \dots, \frac{1}{n_r}\mathbf{J}_{n_r}\right),$$

where $\mathbf{J}_n = \mathbf{1}_n \mathbf{1}_n'$ and $\mathbf{1}_n$ is a $n$-dimensional vector with all the elements being ones. $\qquad \triangle$

# 4  Two sample $t$-test as DOE

We briefly review the two-sample $t$-test statistic. Suppose that $Y_{1j} \overset{iid}{\sim} N(\mu_1, \sigma^2)$ for $j = 1, 2, \ldots, n_1$ and $Y_{2j} \overset{iid}{\sim} N(\mu_2, \sigma^2)$ for $j = 1, 2, \ldots, n_2$. The two-sample $t$-test statistic under $H_0 : \mu_1 = \mu_2$ is given by

$$T = \frac{\overline{Y}_{1\bullet} - \overline{Y}_{2\bullet}}{\sqrt{S_p^2 \cdot \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \sim t(\mathrm{df} = n_1 + n_2 - 2),$$

where

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}, \ S_1^2 = \frac{1}{n_1 - 1}\sum_{j=1}^{n_1}(Y_{1j} - \overline{Y}_{1\bullet})^2, \ S_2^2 = \frac{1}{n_2 - 1}\sum_{j=1}^{n_2}(Y_{2j} - \overline{Y}_{2\bullet})^2.$$

Then $S_p^2$ can be rewritten as

$$S_p^2 = \frac{\displaystyle\sum_{j=1}^{n_1}(Y_{1j} - \overline{Y}_{1\bullet})^2 + \sum_{j=1}^{n_2}(Y_{2j} - \overline{Y}_{2\bullet})^2}{n_1 + n_2 - 2} = \frac{\displaystyle\sum_{i=1}^{r}\sum_{j=1}^{n_i}(Y_{ij} - \overline{Y}_{i\bullet})^2}{N - r} = \frac{\mathrm{SSE}}{N - r} = \mathrm{MSE},$$

where $N = n_1 + n_2$ and $r = 2$. We also have $T^2 \sim F(r - 1, N - r)$ and

$$T^2 = \frac{(\overline{Y}_{1\bullet} - \overline{Y}_{2\bullet})^2}{\mathrm{MSE} \cdot \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} = \frac{(\overline{Y}_{1\bullet} - \overline{Y}_{2\bullet})^2}{\frac{\mathrm{SSE}}{N - r} \cdot \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} = \frac{\frac{n_1 n_2}{N} \cdot (\overline{Y}_{1\bullet} - \overline{Y}_{2\bullet})^2}{\mathrm{SSE}/(N - r)} \sim F(r - 1, N - r). \tag{28}$$

We can also rewrite SStr for $r = 2$ as follows. Substituting $\overline{Y}_{\bullet\bullet} = (n_1 \overline{Y}_{1\bullet} + n_2 \overline{Y}_{2\bullet})/N$ into

$$\mathrm{SStr} = \sum_{i=1}^{r} n_i(\overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet})^2 = n_1(\overline{Y}_{1\bullet} - \overline{Y}_{\bullet\bullet})^2 + n_2(\overline{Y}_{2\bullet} - \overline{Y}_{\bullet\bullet})^2,$$

we have

$$\begin{aligned}
\mathrm{SStr} &= n_1\left(\overline{Y}_{1\bullet} - \frac{n_1\overline{Y}_{1\bullet} + n_2\overline{Y}_{2\bullet}}{N}\right)^2 + n_2\left(\overline{Y}_{2\bullet} - \frac{n_1\overline{Y}_{1\bullet} + n_2\overline{Y}_{2\bullet}}{N}\right)^2 \\
&= n_1\left[\frac{(n_1 + n_2)\overline{Y}_{1\bullet}}{N} - \frac{n_1\overline{Y}_{1\bullet} + n_2\overline{Y}_{2\bullet}}{N}\right]^2 + n_2\left[\frac{(n_1 + n_2)\overline{Y}_{2\bullet}}{N} - \frac{n_1\overline{Y}_{1\bullet} + n_2\overline{Y}_{2\bullet}}{N}\right]^2 \\
&= n_1\left[\frac{n_2(\overline{Y}_{1\bullet} - \overline{Y}_{2\bullet})}{N}\right]^2 + n_2\left[\frac{n_1(\overline{Y}_{2\bullet} - \overline{Y}_{1\bullet})}{N}\right]^2 \\
&= \left(\frac{n_1 n_2^2}{N^2} + \frac{n_2 n_1^2}{N^2}\right)(\overline{Y}_{1\bullet} - \overline{Y}_{2\bullet})^2 = \frac{n_1 n_2}{N}(\overline{Y}_{1\bullet} - \overline{Y}_{2\bullet})^2.
\end{aligned}$$

Comparing the above with (28), we have

$$T^2 = \frac{\mathrm{SStr}/(r - 1)}{\mathrm{SSE}/(N - r)} \sim F(r - 1, N - r),$$

where $r = 2$ again. That is,

$$\boxed{F = \frac{\text{MStr}}{\text{MSE}} \sim F(r-1, N-r).}$$

# 5  Examples

**Example 1.**   As a real-data example, we consider the two-sample $t$-test problem in Example 8.2-2 of Hogg et al. (2015).

> *A machine with filler heads packages a product. This machine has 12 fillers on the left side*
> *and another 12 fillers on the right side. Let $Y_{1j}$ and $Y_{2j}$ be the fill weights in grams when a*
> *machine fills a package by the left and right heads, respectively. We assume that $Y_{1j}$ and $Y_{2j}$*
> *are normally distributed with a common variance. The data sets are given by*
>
> $Y_{1j} = \{1071, 1076, 1070, 1083, 1082, 1067, 1078, 1080, 1075, 1084, 1075, 1080\}$ *and*
>
> $Y_{2j} = \{1074, 1069, 1075, 1067, 1068, 1079, 1082, 1064, 1070, 1073, 1072, 1075\}.$

Ⓡ: Reading Data

```
1  > Y1 = c(1071,1076,1070,1083,1082,1067,1078,1080,1075,1084,1075,1080)
2  > Y2 = c(1074,1069,1075,1067,1068,1079,1082,1064,1070,1073,1072,1075)
3  > boxplot(Y1,Y2, names=c("X","Y") )
```

Ⓡ: Two sample $t$-test

```
1  > t.test(Y1,Y2, alternative="two.sided", var.equal=TRUE)
2          Two Sample t-test
3  data:  Y1 and Y2
4  t = 2.053, df = 22, p-value = 0.05215
5  alternative hypothesis: true difference in means is not equal to 0
6  95 percent confidence interval:
7   -0.04488773  8.87822107
8  sample estimates:
9  mean of x mean of y
10   1076.750  1072.333
```

Ⓡ: ANOVA with `aov()` function with factor

```
1  > group = factor( c(rep("Y1",12), rep("Y2",12)) )
2  > Y = c(Y1,Y2)
3  > plot(Y~group)  # Boxplot
4
5  > myaov = aov(Y ~ group)
6  > summary(myaov)
7             Df Sum Sq Mean Sq F value Pr(>F)
8  group       1  117.0  117.04   4.215 0.0522 .
9  Residuals  22  610.9   27.77
10  ---
11
12  > # ---------------------------------------------
13  > # Parameter Estimation
14  > # ---------------------------------------------
15  > mui = tapply(Y, group, mean)
16  > tau = mui - mean(Y)
17  > cbind(mui, tau)
```

```
18           mui        tau
19  Y1 1076.750   2.208333
20  Y2 1072.333  -2.208333
```

## Ⓡ: ANOVA with `lm()` function with factor

```
1  > group = factor( c(rep("Y1",12), rep("Y2",12)) )
2  > mylm = lm (Y ~ group)
3  > anova(mylm)     # same as the summary(myaov) above
4  Analysis of Variance Table
5  Response: Y
6            Df Sum Sq Mean Sq F value  Pr(>F)
7  group      1 117.04 117.042  4.2148 0.05215 .
8  Residuals 22 610.92  27.769
9  ---
10
11 > summary(mylm)  # ANOVA is OK, but parameter estimates are not
12 Call:
13 lm(formula = Y ~ group)
14 Residuals:
15     Min     1Q  Median     3Q     Max
16 -9.7500 -3.5833  0.1667  3.2500  9.6667
17
18 Coefficients:
19             Estimate Std. Error t value Pr(>|t|)
20 (Intercept) 1076.750      1.521 707.825   <2e-16 ***
21 groupY2        -4.417      2.151  -2.053   0.0522 .
22 ---
23
24 Residual standard error: 5.27 on 22 degrees of freedom
25 Multiple R-squared:  0.1608,    Adjusted R-squared:  0.1226
26 F-statistic: 4.215 on 1 and 22 DF,  p-value: 0.05215
```

## Ⓡ: Regression: cell-means coding

```
1  > z1 =  c(rep(1,12), rep(0,12))
2  > z2 =  c(rep(0,12), rep(1,12))
3  >
4  > # under H0
5  > mylm0 = lm( Y ~ 1 )
6  > anova(mylm0)
7  Analysis of Variance Table
8  Response: Y
9            Df Sum Sq Mean Sq F value Pr(>F)
10 Residuals 23 727.96   31.65
11 >
12 > # under H1
13 > mylm1 = lm( Y ~ 0 + z1 + z2 )
14 > mylm1
15 Call:
16 lm(formula = Y ~ 0 + z1 + z2)
17 Coefficients:
18   z1    z2
19 1077  1072
20
21 > anova(mylm1)
22 Analysis of Variance Table
23 Response: Y
24           Df   Sum Sq  Mean Sq F value     Pr(>F)
25 z1         1 13912687 13912687  501016 < 2.2e-16 ***
26 z2         1 13798785 13798785  496914 < 2.2e-16 ***
27 Residuals 22      611       28
28 ---
29
30 > # ------------------------------------------------
31 > # Parameter Estimation
32 > # ------------------------------------------------
33 > mu0 = mean(Y)
34 > mui = coef(mylm1)
35 > tau = mui - mu0
36 > cbind(mui, tau)
```

```
37          mui        tau
38  z1 1076.750   2.208333
39  z2 1072.333  -2.208333
40  >
41  > # -----------------------------------------------
42  > # Difference between H0 and H1  (Note: anova() function)
43  > # -----------------------------------------------
44  > anova(mylm0, mylm1)
45  Analysis of Variance Table
46  Model 1: Y ~ 1
47  Model 2: Y ~ 0 + z1 + z2
48    Res.Df    RSS Df Sum of Sq      F Pr(>F)
49  1     23 727.96
50  2     22 610.92  1    117.04 4.2148 0.05215 .
51  ---
```

## Ⓡ: Regression: Effect coding

```
1   > z1 =  c(rep(1,12), rep(-1,12))
2   >
3   > # under H0
4   > mylm0 = lm( Y ~ 1 )
5   > anova(mylm0)
6   Analysis of Variance Table
7   Response: Y
8             Df Sum Sq Mean Sq F value Pr(>F)
9   Residuals 23 727.96   31.65
10
11  > # under H1 (note: with intercept)
12  > mylm1 = lm( Y ~ 1 + z1 )
13  > mylm1
14  Call:
15  lm(formula = Y ~ 1 + z1)
16  Coefficients:
17  (Intercept)          z1
18     1074.542       2.208
19
20  > anova(mylm1)
21  Analysis of Variance Table
22  Response: Y
23            Df Sum Sq Mean Sq F value  Pr(>F)
24  z1         1 117.04 117.042  4.2148 0.05215 .
25  Residuals 22 610.92  27.769
26  ---
27
28  > # -----------------------------------------------
29  > # Parameter Estimation
30  > # -----------------------------------------------
31  > mu0 = coef(mylm1)[1]
32  > tau = c( coef(mylm1)[2], -coef(mylm1)[2] )
33  > mui = c(mu0+tau)
34  > cbind(mui, tau)
35          mui        tau
36  z1 1076.750   2.208333
37  z1 1072.333  -2.208333
38
39  > # -----------------------------------------------
40  > # Difference between H0 and H1
41  > # -----------------------------------------------
42  > anova(mylm0, mylm1)
43  Analysis of Variance Table
44  Model 1: Y ~ 1
45  Model 2: Y ~ 1 + z1
46    Res.Df    RSS Df Sum of Sq      F Pr(>F)
47  1     23 727.96
48  2     22 610.92  1    117.04 4.2148 0.05215 .
49  ---
```

**Example 2.**   As a real-data example, we consider Example 3.6 of the textbook by Kim (2014).

> *To determine the effect of iris color on critical flicker frequency (CFF), iris colours and CFFs*
>
> *for 19 persons are recorded. The subjects are divided into three groups on the basis of iris color*
>
> *(blue, brown and green). Let $Y_{1j}$, $Y_{2j}$ and $Y_{3j}$ be the CFFs of blue, brown and green iris colors,*
>
> *respectively. We assume that $Y_{1j}$, $Y_{2j}$ and $Y_{3j}$ are normally distributed with a common variance.*
>
> *As will be shown below, the analysis indicates that iris color is a statistically significant factor*
>
> *for CFF. The original data set is provided in Smith and Misiak (1973).*

Ⓡ: Reading Data

```
1  > Y1 = c(25.7, 27.2, 29.9, 28.5, 29.4, 28.3)          # Blue
2  > Y2 = c(26.8, 27.9, 23.7, 25, 26.3, 24.8, 25.7, 24.5) # Brown
3  > Y3 = c(26.4, 24.2, 28.0, 26.9, 29.1)                # Green
4  > Y = c(Y1, Y2, Y3)
5  > n1 = length(Y1); n2 = length(Y2); n3 = length(Y3)   # no. of subjects
```

Ⓡ: ANOVA with `aov()` function with factor

```
1  > # Check the input
2  > color = factor( rep(c("Blue","Brown","Green"), c(n1,n2,n3)) )
3  > levels(color)
4  [1] "Blue"  "Brown" "Green"
5  > cbind(Y,color)
6          Y color
7   [1,] 25.7     1
8   [2,] 27.2     1
9        .........
10       .........
11 # Plots
12 # stripchart(Y ~ color, vertical=T)
13 # boxplot(Y ~ color, ylab="Flicker")
14
15 > myaov = aov(Y ~ color)
16 > summary(myaov)
17            Df Sum Sq Mean Sq F value Pr(>F)
18 color       2  23.00  11.499   4.802 0.0232 *
19 Residuals  16  38.31   2.394
20 ---
21
22 > # mu.hat and tau.hat
23 > mui = tapply(Y, color, mean)
24 > tau = mui - mean(Y)
25 > cbind(mui, tau)
26          mui        tau
27 Blue  28.16667  1.4140351
28 Brown 25.58750 -1.1651316
29 Green 26.92000  0.1673684
30
31 # Diagnostic
32 # par ( mfrow=c(2,2) )
33 # plot(myaov)
```

Ⓡ: ANOVA with `lm()` function with factor

```
1  > color = factor( rep(c("Blue","Brown","Green"), c(n1,n2,n3)) )
2  > mylm = lm (Y ~ color)
3  > anova(mylm)    # same as the summary(myaov) above
4  Analysis of Variance Table
```

```
5   Response: Y
6             Df Sum Sq Mean Sq F value  Pr(>F)
7   color      2 22.997 11.4986  4.8023 0.02325 *
8   Residuals 16 38.310  2.3944
9   ---
10
11  > summary(mylm)  # ANOVA is OK, but parameter estimates are not
12  Call:
13  lm(formula = Y ~ color)
14  Residuals:
15      Min     1Q  Median      3Q     Max
16  -2.7200 -0.8771  0.1125  1.1462  2.3125
17
18  Coefficients:
19              Estimate Std. Error t value Pr(>|t|)
20  (Intercept) 28.1667     0.6317  44.588  < 2e-16 ***
21  colorBrown  -2.5792     0.8357  -3.086  0.00708 **
22  colorGreen  -1.2467     0.9370  -1.331  0.20200
23  ---
24
25  Residual standard error: 1.547 on 16 degrees of freedom
26  Multiple R-squared:  0.3751,    Adjusted R-squared:  0.297
27  F-statistic: 4.802 on 2 and 16 DF,  p-value: 0.02325
```

## Ⓡ: Regression: cell-means coding

```
1   > z1 = rep( c(1,0,0), c(n1,n2,n3) )
2   > z2 = rep( c(0,1,0), c(n1,n2,n3) )
3   > z3 = rep( c(0,0,1), c(n1,n2,n3) )
4
5   > # under H0
6   > mylm0 = lm( Y ~ 1 )
7   > anova(mylm0)
8   Analysis of Variance Table
9   Response: Y
10            Df Sum Sq Mean Sq F value Pr(>F)
11  Residuals 18 61.307   3.406
12
13  > # under H1 (Note: no-intercept model)
14  > mylm1 = lm( Y ~ 0 + z1 + z2 + z3 )
15  > anova(mylm1)
16  Analysis of Variance Table
17  Response: Y
18            Df Sum Sq Mean Sq F value    Pr(>F)
19  z1         1 4760.2  4760.2  1988.1 < 2.2e-16 ***
20  z2         1 5237.8  5237.8  2187.5 < 2.2e-16 ***
21  z3         1 3623.4  3623.4  1513.3 < 2.2e-16 ***
22  Residuals 16   38.3     2.4
23  ---
24
25  > # mu.hat and tau.hat
26  > mui = coef(mylm1)
27  > tau = mui - mean(Y)
28  > cbind(mui, tau)
29         mui        tau
30  z1 28.16667  1.4140351
31  z2 25.58750 -1.1651316
32  z3 26.92000  0.1673684
33
34  > # Difference between H0 and H1  (Note: anova() function)
35  > anova(mylm0, mylm1)
36  Analysis of Variance Table
37  Model 1: Y ~ 1
38  Model 2: Y ~ 0 + z1 + z2 + z3
39    Res.Df    RSS Df Sum of Sq      F  Pr(>F)
40  1     18 61.307
41  2     16 38.310  2    22.997 4.8023 0.02325 *
42  ---
```

## Ⓡ: Regression: effect coding

```
1   > # Wrong Effect coding (factor-effect model)
2   > z1 = rep( c(1,0,-1), c(n1,n2,n3) ) # This is only for balanced sample
3   > z2 = rep( c(0,1,-1), c(n1,n2,n3) ) # This is only for balanced sample
4
5   > # Correct Effect coding
6   > z1 = rep( c(1,0,-n1/n3), c(n1,n2,n3) )
7   > z2 = rep( c(0,1,-n2/n3), c(n1,n2,n3) )
8
9   > # under H0
10  > mylm0 = lm( Y ~ 1 )
11  > anova(mylm0)
12  Analysis of Variance Table
13  Response: Y
14           Df Sum Sq Mean Sq F value Pr(>F)
15  Residuals 18 61.307   3.406
16
17  > # under H1 (note: with intercept)
18  > mylm1 = lm( Y ~ 1 + z1 + z2)
19  > mylm1   # beta0 is a grand mean.
20  Call:
21  lm(formula = Y ~ 1 + z1 + z2)
22  Coefficients:
23  (Intercept)           z1           z2
24      26.753        1.414      -1.165
25
26  > anova(mylm1)
27  Analysis of Variance Table
28  Response: Y
29           Df Sum Sq Mean Sq F value  Pr(>F)
30  z1        1  4.239  4.2387  1.7703 0.20200
31  z2        1 18.759 18.7586  7.8344 0.01287 *
32  Residuals 16 38.310  2.3944
33  ---
34
35  > mu0 = coef(mylm1)[1]
36  > c(mu0, mean(Y))
37  (Intercept)
38     26.75263    26.75263
39  > tau = c( coef(mylm1)[-1], -sum(c(n1,n2)*coef(mylm1)[-1])/n3 )
40  > mui = c(mu0+tau)
41  > cbind(mui, tau)
42          mui        tau
43  z1 28.16667  1.4140351
44  z2 25.58750 -1.1651316
45     26.92000  0.1673684
46
47  > # Difference between H0 and H1
48  > # Note: ANOVA decomposition is the same as the cell-means coding
49  > anova(mylm0, mylm1)
50  Analysis of Variance Table
51  Model 1: Y ~ 1
52  Model 2: Y ~ 1 + z1 + z2
53    Res.Df    RSS Df Sum of Sq      F  Pr(>F)
54  1     18 61.307
55  2     16 38.310  2    22.997 4.8023 0.02325 *
```

$\|$

**Example 3.** Revisit Example 2. In Section 1.2, we considered the confidence interval of $\mu_i$. One may obtain the confidence interval using `t.test()` function in R language. But, this is *not* recommended because it does not use all the samples to estimate $\sigma^2$.

Again, the interval estimation of $\mu_i$ with $100(1-\alpha)\%$ coverage is given by

$$\overline{Y}_{i\bullet} \;\pm\; t(1-\tfrac{\alpha}{2};N-r)\sqrt{\frac{\text{MSE}}{n_i}},$$

Note that the degrees of freedom is $N-r$ when all the samples are used. On the other hand, the interval estimation of $\mu_i$ using the $i$th sample only is give by

$$\overline{Y}_{i\bullet} \;\pm\; t(1-\tfrac{\alpha}{2};n_i-1)\sqrt{\frac{S_i^2}{n_i}},$$

As an illustration, we obtain the 95% confidence interval of $\mu_1$ as below.

### Ⓡ: Reading Data

```
1  > Y1 = c(25.7, 27.2, 29.9, 28.5, 29.4, 28.3)        # Blue
2  > Y2 = c(26.8, 27.9, 23.7, 25, 26.3, 24.8, 25.7, 24.5) # Brown
3  > Y3 = c(26.4, 24.2, 28.0, 26.9, 29.1)              # Green
4  > Y = c(Y1, Y2, Y3)
5  > n1 = length(Y1); n2 = length(Y2); n3 = length(Y3)   # no. of subjects
```

### Ⓡ: CI from t.test()

```
1  > # Bad CI
2  > t.test(Y1)
3          One Sample t-test
4  data:  Y1
5  t = 45.154, df = 5, p-value = 1.006e-07
6  alternative hypothesis: true mean is not equal to 0
7  95 percent confidence interval:
8   26.56317 29.77016
9  sample estimates:
10 mean of x
11  28.16667
```

### Ⓡ: CI using $S_1$ only.

```
1  > # Same as CI from t.test()
2  > a = 0.05
3  > D = qt(1-a/2,df=n1-1) * sqrt(var(Y1)/n1)
4  > c(mean(Y1)-D, mean(Y1)+D)
5  [1] 26.56317 29.77016
```

### Ⓡ: CI using MSE (better).

```
1  > # Better CI
2  > a = 0.05; df = n1+n2+n3-3
3  > SSE = sum((Y1-mean(Y1))^2) + sum((Y2-mean(Y2))^2) + sum((Y3-mean(Y3))^2)
4  > MSE = SSE / df
5  > D = qt(1-a/2,df=df) * sqrt(MSE/n1)
6  > c(mean(Y1)-D, mean(Y1)+D)
7  [1] 26.82749 29.50584
```

Again, the $100(1-\alpha)\%$ confidence interval of $\mu_\ell - \mu_m$ is given by

$$\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet} \;\pm\; t(1-\tfrac{\alpha}{2};N-r)\sqrt{\text{MSE}\left(\frac{1}{n_\ell} + \frac{1}{n_m}\right)}.$$

On the other hand, based on the two-sample $t$-test statistic, the confidence interval of $\mu_\ell - \mu_m$ is given by

$$\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet} \ \pm \ t(1 - \tfrac{\alpha}{2}; n_\ell + n_m - 2)\sqrt{S_p^2 \left( \frac{1}{n_\ell} + \frac{1}{n_m} \right)},$$

which, however, is not recommended because it $S_p^2$ does not incorporate all the samples provided. As an illustration, we obtain the 95% confidence interval of $\mu_1 - \mu_2$ as below.

Ⓡ: CI from `t.test()`

```
1  > # Bad CI
2  > t.test(Y1,Y2, var.equal=TRUE)
3
4          Two Sample t-test
5  data:  Y1 and Y2
6  t = 3.3272, df = 12, p-value = 0.006029
7  alternative hypothesis: true difference in means is not equal to 0
8  95 percent confidence interval:
9   0.8902227 4.2681106
10 sample estimates:
11 mean of x mean of y
12  28.16667  25.58750
```

Ⓡ: CI using $S_p^2$ (with two samples).

```
1  > # Same as CI from t.test()
2  > a = 0.05
3  > Sp2 = ((n1-1)*var(Y1)+(n2-1)*var(Y2))/(n1+n2-2)
4  > D = qt(1-a/2,df=n1+n2-2) * sqrt(Sp2*(1/n1+1/n2))
5  > c(mean(Y1)-mean(Y2)-D, mean(Y1)-mean(Y2)+D)
6  [1] 0.8902227 4.2681106
```

Ⓡ: CI using MSE (better).

```
1  > a = 0.05; df = n1+n2+n3-3
2  > SSE = sum((Y1-mean(Y1))^2) + sum((Y2-mean(Y2))^2) + sum((Y3-mean(Y3))^2)
3  > MSE = SSE / df
4  > D = qt(1-a/2,df=df) * sqrt(MSE*(1/n1+1/n2))
5  > c(mean(Y1)-mean(Y2)-D, mean(Y1)-mean(Y2)+D)
6  [1] 0.8076044 4.3507289
```

$\parallel$

# 6 Comparisons among treatment means

## 6.1 Contrasts

When we compare the treatment means, the idea of a contrast is widely used because we can compare the treatment means using them. As an illustration, suppose that one would like to test the hypothesis

$$H_0 : \mu_1 = \mu_2 \ \ \text{versus} \ \ H_1 : \mu_1 \neq \mu_2,$$

which is equivalent to

$$H_0 : \mu_1 - \mu_2 = 0 \quad \text{versus} \quad H_1 : \mu_1 - \mu_2 \neq 0.$$

Thus, this hypothesis testing can be carried out by using a linear combination of the parameters. This linear combination is called a contrast.

---

**Definition 2.** We denote $\boldsymbol{\theta} = (\mu_1, \mu_2, \ldots, \mu_r)$ be a set of parameters (or statistics) and $\mathbf{a} = (a_1, a_2, \ldots, a_r)$ be a collection of known constants with its sum begin zero, that is,

$$\sum_{i=1}^{r} a_i = 0.$$

Then $\mathbf{a} = (a_1, a_2, \ldots, a_r)$ is called *contrast constants* and its linear combination below is called a *contrast*

$$\Gamma = \sum_{i=1}^{r} a_i \mu_i.$$

---

For example, with $\mathbf{a} = (1, -1, 0, \ldots, 0)$, we have

$$\Gamma = \sum_{i=1}^{r} a_i \mu_i = \mu_1 - \mu_2.$$

Thus, the hypothesis testing $H_0 : \mu_1 = \mu_2$ versus $H_1 : \mu_1 \neq \mu_2$ can be rewritten as a contrast as below

$$H_0 : \Gamma = 0 \quad \text{versus} \quad H_1 : \Gamma \neq 0.$$

---

**Definition 3.** Let $\overline{Y}_{i\bullet}$ be the estimator of the $i$ treatment $\mu_i$ based on $n_i$ observations and $\mathbf{a} = (a_1, a_2, \ldots, a_r)$ and $\mathbf{b} = (b_1, b_2, \ldots, b_r)$ be contrast constants satisfying

$$\sum_{i=1}^{r} \frac{a_i b_i}{n_i} = 0.$$

Then two contrasts below are called *orthogonal*

$$\Gamma_{\mathbf{a}} = \sum_{i=1}^{r} a_i \mu_i \quad \text{and} \quad \Gamma_{\mathbf{b}} = \sum_{i=1}^{r} b_i \mu_i.$$

---

**Remark 7.** Some textbooks state that $\Gamma_{\mathbf{a}}$ and $\Gamma_{\mathbf{b}}$ are *orthogonal* especially when $\sum_{i=1}^{r} a_i b_i = 0$ and *uncorrelated* when $\sum_{i=1}^{r} a_i b_i / n_i = 0$. Of course, if the balanced design is used (that is, $n_1 = n_2 = \cdots = n_r = n$), then the uncorrelation condition becomes $\sum_{i=1}^{r} a_i b_i = 0$.          $\triangle$

## 6.2  Inferences regarding contrasts

We assumed that $Y_{ij} = \mu_i + \epsilon_{ij}$, where $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, $i = 1, 2, \ldots, r$, and $j = 1, 2, \ldots, n_i$ so that

$$\overline{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \sim N\left(\mu_i, \frac{\sigma^2}{n_i}\right). \tag{29}$$

Thus, for any contrast constants $\mathbf{a} = (a_1, a_2, \ldots, a_r)$, $\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet}$ is also normally distributed with

$$E\left(\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet}\right) = \sum_{i=1}^{r} a_i \mu_i \ \ \text{and} \ \ \text{Var}\left(\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet}\right) = \sigma^2 \sum_{i=1}^{r} \frac{a_i^2}{n_i}.$$

Standardizing $\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet}$, we have

$$\frac{\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet} - \sum_{i=1}^{r} a_i \mu_i}{\sqrt{\sigma^2 \sum_{i=1}^{r} a_i^2 / n_i}} \sim N(0, 1). \tag{30}$$

Note that (30) defines a pivot but it includes a nuisance parameter $\sigma^2$.

---

**Lemma 4.**  *Under the assumption that $Y_{ij} = \mu_i + \epsilon_{ij}$ with $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, $\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet} - \sum_{i=1}^{r} a_i \mu_i$ and SSE in (19) are independent.*

---

*Proof.*  The proof is very similar to that of Lemma 2 (c).  □

Recall that $\text{SSE}/\sigma^2 \sim \chi^2(N - r)$ from (19). Thus, using Lemma 4, we can Studentize (30) as

$$\frac{\frac{\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet} - \sum_{i=1}^{r} a_i \mu_i}{\sqrt{\sigma^2 \sum_{i=1}^{r} a_i^2 / n_i}}}{\sqrt{\frac{1}{\sigma^2} \cdot \text{SSE}/(N - r)}} = \frac{\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet} - \sum_{i=1}^{r} a_i \mu_i}{\sqrt{\text{MSE} \sum_{i=1}^{r} a_i^2 / n_i}} \sim t(\text{df} = N - r). \tag{31}$$

Note that MSE can be viewed as an extended version of the pooled sample variance because

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2 + \cdots + (n_r - 1)S_r^2}{(n_1 - 1) + (n_2 - 1) + \cdots + (n_r - 1)} = \frac{\text{SSE}}{N - r},$$

where $S_i^2 = \sum_{j=1}^{n_i} (Y_{ij} - \overline{Y}_{i\bullet})^2 / (n_i - 1)$ is the sample variance of the $i$th sample. One can perform the hypothesis testing of

$$H_0 : \sum_{i=1}^{r} a_i \mu_i = 0 \text{ versus } H_1 : \sum_{i=1}^{r} a_i \mu_i \neq 0$$

by rejecting $H_0$ if

$$\left| \frac{\sum_{i=1}^{r} a_i \overline{Y_{i\bullet}}}{\sqrt{\text{MSE} \sum_{i=1}^{r} a_i^2/n_i}} \right| > t(1 - \tfrac{\alpha}{2}; N - r)$$

at the significance level of $\alpha$.

Both (30) and (31) define a pivot, but the former includes a nuisance parameter. Thus, we invert (31) for the contrast $\sum_{i=1}^{r} a_i \mu_i$ to obtain an interval estimator whose endpoints with $100(1-\alpha)\%$ confidence level are given by

$$\sum_{i=1}^{r} a_i \overline{Y_{i\bullet}} \ \pm \ t(1 - \tfrac{\alpha}{2}; N - r) \cdot \sqrt{\text{MSE} \sum_{i=1}^{r} \frac{a_i^2}{n_i}}. \tag{32}$$

It should be noted that the endpoints in (14) are easily obtained from (32) with $a_\ell = 1$, $a_m = -1$ and $a_i = 0$ for $i \neq \ell, m$.

---

**Lemma 5.**   *Under the assumption that $Y_{ij} = \mu_i + \epsilon_{ij}$ with $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, we have*

$$\text{Cov}\left( \sum_{i=1}^{r} a_i \overline{Y_{i\bullet}}, \sum_{j=1}^{r} b_j \overline{Y_{j\bullet}} \right) = \sigma^2 \sum_{i=1}^{r} \frac{a_i b_i}{n_i}.$$

---

*Proof.*   It is immediate from Lemma 1 that

$$\text{Cov}\left( \sum_{i=1}^{r} a_i \overline{Y_{i\bullet}}, \sum_{j=1}^{r} b_j \overline{Y_{j\bullet}} \right) = \sum_{i=1}^{r} \sum_{j=1}^{r} a_i b_j \text{Cov}\left( \overline{Y_{i\bullet}}, \overline{Y_{j\bullet}} \right)$$

Since $\overline{Y_{i\bullet}}$ and $\overline{Y_{j\bullet}}$ are independent for $i \neq j$, the above becomes

$$\sum_{i=1}^{r} a_i b_i \text{Cov}\left( \overline{Y_{i\bullet}}, \overline{Y_{i\bullet}} \right) = \sum_{i=1}^{r} a_i b_i \text{Var}(\overline{Y_{i\bullet}}) = \sum_{i=1}^{r} a_i b_i \frac{\sigma^2}{n_i},$$

which completes the proof.                                                                                       $\square$

Hence, the contrasts are uncorrelated if they are orthogonal by Definition 3.

---

**Definition 4.** For a contrast $\Gamma = \sum_{i=1}^{r} a_i \mu_i$ with the estimate $\hat{\Gamma} = \sum_{i=1}^{r} a_i \overline{Y_{i\bullet}}$, the term

$$\frac{\left( \sum_{i=1}^{r} a_i \overline{Y_{i\bullet}} \right)^2}{\sum_{i=1}^{r} a_i^2/n_i}$$

is called the *contrast sum of squares* due to the contrast $\hat{\Gamma}$.

---

Note that it is immediate from (30) that under $H_0 : \sum_{i=1}^{r} a_i \mu_i = 0$

$$\frac{1}{\sigma^2} \frac{\left(\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} a_i^2/n_i} \sim \chi^2(1).$$

In general, if the sum of squares has a chi-squared distribution with $\nu$ degrees of freedom, we can decompose it into $\nu$ independent chi-squared random variables with all one degree of freedom using orthogonal (uncorrelated) contrasts. Thus, the treatment sum of squares in (15) can be decomposed into $r-1$ contrast sums of squares with uncorrelated contrasts.

If there are $r$ treatments, we can find $r-1$ orthogonal (uncorrelated) sets of contrast constants. Denote these sets by $\mathbf{a}^{(\ell)} = (a_1^{(\ell)}, a_2^{(\ell)}, \ldots, a_r^{(\ell)})$ for $\ell = 1, 2, \ldots, r-1$ which satisfy Definition 3, that is,

$$\sum_{i=1}^{r} \frac{a_i^{(\ell)} a_i^{(\ell')}}{n_i} = 0,$$

for all $\ell \neq \ell'$. Note that $\left[\mathbf{1}_r, \mathbf{a}^{(1)'}, \mathbf{a}^{(2)'}, \ldots, \mathbf{a}^{(r-1)'}\right]$ constitute a $r \times r$ square matrix of full rank $r$, where $\mathbf{1}_r$ is a $r$-dimensional vector with all the elements being ones and we also have $\mathbf{1}_r' \mathbf{a}^{(\ell)} = 0$ for $\ell = 1, 2, \ldots, r-1$. Then we can decompose SStr in (15) into contrast sums of squares as follows

$$\sum_{i=1}^{r} n_i (\overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet})^2 = \frac{\left(\sum_{i=1}^{r} a_i^{(1)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(1)}\right)^2/n_i} + \frac{\left(\sum_{i=1}^{r} a_i^{(2)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(2)}\right)^2/n_i} + \cdots + \frac{\left(\sum_{i=1}^{r} a_i^{(r-1)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(r-1)}\right)^2/n_i}$$

$$= \sum_{\ell=1}^{r-1} \frac{\left(\sum_{i=1}^{r} a_i^{(\ell)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(\ell)}\right)^2/n_i}. \tag{33}$$

**Example 4.** As a real-data example, we consider Example 3.4 of the textbook by Kim (2014).

Ⓡ: Reading Data

```
1  > Y1 = c(2.4, 2.7, 3.1, 3.1)
2  > Y2 = c(0.7, 1.6, 1.7, 1.8)
3  > Y3 = c(2.4, 3.1, 5.4, 6.1)
4  > Y4 = c(0.3, 0.3, 2.4, 2.7)
5  > Y5 = c(0.5, 0.9, 1.4, 2.0)
6  > Y = c(Y1, Y2, Y3, Y4, Y5)
7  > n = length(Y1)
8  > treat = factor( rep(LETTERS[1:5], rep(n,5)) )
```

Ⓡ: Calculating SStr

```
1  # Calculating SStr with aov() function
2  > myaov = aov(Y~treat)
3  > summary(myaov)
```

```
 4                 Df Sum Sq Mean Sq F value  Pr(>F)
 5  treat          4  27.01   6.752    5.966 0.00444 **
 6  Residuals     15  16.98   1.132
 7  ---
 8  > SStr = summary(myaov)[[1]][["Sum Sq"]][1]
 9  > SStr
10  [1] 27.007
11
12  # Calculating SStr manually
13  > Yi.bar = c(mean(Y1), mean(Y2), mean(Y3), mean(Y4), mean(Y5))
14  > Ybarbar = mean(Y)
15  > SStr = sum( n*(Yi.bar-Ybarbar)^2 )
16  > SStr
17  [1] 27.007
```

## Ⓡ: Decomposition of SStr with Contrasts

```
 1  > # ---------------------------------------------
 2  > # Orthogonal contrasts (manual contrasts)
 3  > # ---------------------------------------------
 4  > a1 = c(1, -1/4, -1/4, -1/4, -1/4)
 5  > a2 = c(0,    1, -1/3, -1/3, -1/3)
 6  > a3 = c(0,    0,    1, -1/2, -1/2)
 7  > a4 = c(0,    0,    0,    1,   -1)
 8  > contrasts(treat) = cbind(a1,a2,a3,a4)
 9  > contrasts(treat)
10       a1         a2    a3 a4
11  A  1.00  0.0000000  0.0  0
12  B -0.25  1.0000000  0.0  0
13  C -0.25 -0.3333333  1.0  0
14  D -0.25 -0.3333333 -0.5  1
15  E -0.25 -0.3333333 -0.5 -1
16
17  > # Check the validity of contrasts (sum=0)
18  > cnt = contrasts(treat)
19  > apply(cnt, 2, sum)
20            a1            a2            a3            a4
21  0.000000e+00 5.551115e-17 0.000000e+00 0.000000e+00
22
23  > # Check the orthogonality
24  > crossprod(cnt)
25                a1            a2  a3 a4
26  a1  1.250000e+00 -2.775558e-17 0.0  0
27  a2 -2.775558e-17  1.333333e+00 0.0  0
28  a3  0.000000e+00  0.000000e+00 1.5  0
29  a4  0.000000e+00  0.000000e+00 0.0  2
30
31  > # Check the decomposition of SStr
32  > myaov = aov( Y ~ treat )
33  > summary.lm(myaov)
34  Call:
35  aov(formula = Y ~ treat)
36  Residuals:
37      Min     1Q  Median      3Q     Max
38  -1.8500 -0.7125  0.1750  0.4625  1.8500
39  Coefficients:
40             Estimate Std. Error t value Pr(>|t|)
41  (Intercept)   2.2300     0.2379   9.375 1.16e-07 ***
42  treata1       0.5950     0.4757   1.251 0.230213
43  treata2      -0.6313     0.4606  -1.370 0.190729
44  treata3       1.9583     0.4343   4.509 0.000416 ***
45  treata4       0.1125     0.3761   0.299 0.768957
46  ---
47  Residual standard error: 1.064 on 15 degrees of freedom
48  Multiple R-squared:  0.614,     Adjusted R-squared:  0.5111
49  F-statistic: 5.966 on 4 and 15 DF,  p-value: 0.004442
50
51  > MSE = summary(myaov)[[1]][["Mean Sq"]][2]
52  > MSE
53  [1] 1.131667
54  > t.stat = summary.lm(myaov)[[4]][-1,3]
```

```
55  > SStr.for.contrast = t.stat^2 * MSE
56  > SStr.for.contrast
57    treata1    treata2    treata3    treata4
58   1.770125   2.125208  23.010417   0.101250
59  > sum(SStr.for.contrast)  # Note: SStr = 27.007
60  [1] 27.007
```

We have used the following orthogonal contrast constants: $\mathbf{a}^{(1)} = (1, -\frac{1}{4}, -\frac{1}{4}, -\frac{1}{4}, -\frac{1}{4})$, $\mathbf{a}^{(2)} = (0, 1, -\frac{1}{3}, -\frac{1}{3}, -\frac{1}{3})$, $\mathbf{a}^{(3)} = (0, 0, 1, -\frac{1}{2}, -\frac{1}{2})$, and $\mathbf{a}^{(4)} = (0, 0, 0, 1, -1)$. As shown in the R program above, the treatment sum of squares was decomposed as follows

$$27.007 \ (\text{SStr}) = 1.770125 + 2.125208 + 23.010417 + 0.101250.$$

The R program reports only $t$-test statistic

$$T_\ell = \frac{\sum_{i=1}^{r} a_i^{(\ell)} \overline{Y}_{i\bullet}}{\sqrt{\text{MSE} \sum_{i=1}^{r} \left(a_i^{(\ell)}\right)^2 / n_i}} \tag{34}$$

which is from (31) under $H_0 : \sum_{i=1}^{r} a_i^{(\ell)} \mu_i = 0$. Along with the above orthogonal contrast constants, we can test the below

$T_1$ for testing $H_0^{(1)} : \ \mu_1 = (\mu_2 + \mu_3 + \mu_4 + \mu_5)/4,$

$T_2$ for testing $H_0^{(2)} : \ \mu_2 = (\mu_3 + \mu_4 + \mu_5)/3,$

$T_3$ for testing $H_0^{(3)} : \ \mu_3 = (\mu_4 + \mu_5)/2,$

$T_4$ for testing $H_0^{(4)} : \ \mu_4 = \mu_5.$

Ⓡ: Decomposition of SStr with Contrasts (using $t$-test statistics)

```
1   > # Check the t-test statistics
2   > T1 = sum(a1*Yi.bar) / sqrt(MSE* sum(a1^2/n) )
3   > T2 = sum(a2*Yi.bar) / sqrt(MSE* sum(a2^2/n) )
4   > T3 = sum(a3*Yi.bar) / sqrt(MSE* sum(a3^2/n) )
5   > T4 = sum(a4*Yi.bar) / sqrt(MSE* sum(a4^2/n) )
6   > c(T1,T2,T3,T4)
7   [1]  1.250670 -1.370382  4.509236  0.299115
8   > c(T1,T2,T3,T4)^2 * MSE
9   [1]  1.770125  2.125208 23.010417  0.101250
10  > sum( c(T1,T2,T3,T4)^2 * MSE )
11  [1] 27.007
```

We calculated the $t$-test statistics manually as above, but the `summary.lm()` function reported them as shown earlier. Since it is easily seen from (34) that

$$\frac{\left(\sum_{i=1}^{r} a_i \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} a_i^2 / n_i} = T^2 \cdot \text{MSE},$$

we can calculate the contrast sum of squares using the $t$-test statistic and decompose SStr

as below.

$$\frac{\left(\sum_{i=1}^{r} a_i^{(1)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(1)}\right)^2 / n_i} = T_1^2 \cdot \text{MSE} = 1.250670^2 \times 1.131667 = 1.770125$$

$$\frac{\left(\sum_{i=1}^{r} a_i^{(2)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(2)}\right)^2 / n_i} = T_2^2 \cdot \text{MSE} = (-1.370382)^2 \times 1.131667 = 2.125208$$

$$\frac{\left(\sum_{i=1}^{r} a_i^{(3)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(3)}\right)^2 / n_i} = T_3^2 \cdot \text{MSE} = 4.509236^2 \times 1.131667 = 23.010417$$

$$\frac{\left(\sum_{i=1}^{r} a_i^{(4)} \overline{Y}_{i\bullet}\right)^2}{\sum_{i=1}^{r} \left(a_i^{(4)}\right)^2 / n_i} = T_4^2 \cdot \text{MSE} = 0.299115^2 \times 1.131667 = 0.101250$$

$\parallel$

Contrasts are not uniquely determined. Thus, they should be determined considering the statistical hypothesis of interest. There are several well-known contrasts such as Helmert contrasts (`contr.helmert()`), polynomial contrasts (`contr.poly()`), sum-to-zero contrasts (`contr.sum()`), etc. The sum-to-zero contrasts are similar to effect coding in regression so they are also called effect-coding contrasts. The polynomial contrasts are widely used for testing polynomial patterns in treatment means with more than two treatments (e.g., linear, quadratic, cubic, quartic, etc.). The Helmert and polynomial contrasts are orthogonal contrasts while sum-to-zero contrasts are not. Note that R language also provides the other functions for contrasts such as `contr.treatment()` and `contr.SAS()`, but these are not perpendicular to the intercept vector $\mathbf{1}_r$ which implies that $\mathbf{1}_r' \mathbf{a} = \sum_{i=1}^{r} a_i \neq 0$. Thus, `contr.treatment()` and `contr.SAS()` can not be considered as contrasts in this course!

**Example 5.** Revisit Example 4. In this example, we analyze the data again with the Helmert, polynomial and sum-to-zero contrasts.

Ⓡ: Reading Data

```
1  > Y1 = c(2.4, 2.7, 3.1, 3.1)
2  > Y2 = c(0.7, 1.6, 1.7, 1.8)
3  > Y3 = c(2.4, 3.1, 5.4, 6.1)
4  > Y4 = c(0.3, 0.3, 2.4, 2.7)
5  > Y5 = c(0.5, 0.9, 1.4, 2.0)
6  > Y = c(Y1, Y2, Y3, Y4, Y5)
7  > n = length(Y1)
8  > treat = factor( rep(LETTERS[1:5], rep(n,5)) )
```

Ⓡ: (1) contrasts using `contr.helmert()` function

```
1  > contrasts(treat) = contr.helmert(5)
2  > contrasts(treat)
3    [,1] [,2] [,3] [,4]
4  A   -1   -1   -1   -1
```

```
 5  B     1    -1    -1    -1
 6  C     0     2    -1    -1
 7  D     0     0     3    -1
 8  E     0     0     0     4
 9  >
10  > # Check the validity of contrasts (sum=0)
11  > cnt = contrasts(treat)
12  > apply(cnt, 2, sum)
13  [1] 0 0 0 0
14
15
16  > # Check the orthogonality
17  > crossprod(cnt)
18       [,1] [,2] [,3] [,4]
19  [1,]    2    0    0    0
20  [2,]    0    6    0    0
21  [3,]    0    0   12    0
22  [4,]    0    0    0   20
23
24  > # Check the decomposition of SStr
25  > myaov = aov( Y ~ treat )
26  > summary.lm(myaov)
27  Call:
28  aov(formula = Y ~ treat)
29  Residuals:
30      Min     1Q  Median      3Q     Max
31  -1.8500 -0.7125  0.1750  0.4625  1.8500
32  Coefficients:
33              Estimate Std. Error t value Pr(>|t|)
34  (Intercept)   2.2300     0.2379   9.375 1.16e-07 ***
35  treat1       -0.6875     0.3761  -1.828  0.08752 .
36  treat2        0.7042     0.2171   3.243  0.00546 **
37  treat3       -0.3542     0.1535  -2.307  0.03577 *
38  treat4       -0.2575     0.1189  -2.165  0.04692 *
39  ---
40  Residual standard error: 1.064 on 15 degrees of freedom
41  Multiple R-squared:  0.614,      Adjusted R-squared:  0.5111
42  F-statistic: 5.966 on 4 and 15 DF,  p-value: 0.004442
43
44  > MSE = summary(myaov)[[1]][["Mean Sq"]][2]
45  > t.stat = summary.lm(myaov)[[4]][-1,3]
46  > SStr.for.contrast = t.stat^2 * MSE
47  > SStr.for.contrast
48     treat1     treat2     treat3     treat4
49   3.781250 11.900417   6.020833   5.304500
50  > sum(SStr.for.contrast)  ## Note: SStr = 27.007
51  [1] 27.007
```

### Ⓡ: (2) contrasts using `contr.poly()` function

```
 1  > contrasts(treat) = contr.poly(5)
 2  > contrasts(treat)
 3             .L          .Q            .C            ^4
 4  A -0.6324555  0.5345225 -3.162278e-01  0.1195229
 5  B -0.3162278 -0.2672612  6.324555e-01 -0.4780914
 6  C  0.0000000 -0.5345225 -4.095972e-16  0.7171372
 7  D  0.3162278 -0.2672612 -6.324555e-01 -0.4780914
 8  E  0.6324555  0.5345225  3.162278e-01  0.1195229
 9
10  > # Check the validity of contrasts (sum=0)
11  > cnt = contrasts(treat)
12  > apply(cnt, 2, sum)
13            .L           .Q           .C           ^4
14  0.000000e+00 1.110223e-16 9.001589e-17 6.938894e-17
15
16  > # Check the orthogonality
17  > crossprod(cnt)
18               .L           .Q           .C           ^4
19  .L  1.000000e+00 -1.110223e-16  5.551115e-17 -2.081668e-16
20  .Q -1.110223e-16  1.000000e+00  8.326673e-17 -1.665335e-16
21  .C  5.551115e-17  8.326673e-17  1.000000e+00  1.387779e-17
```

```
22  ^4 -2.081668e-16 -1.665335e-16 1.387779e-17  1.000000e+00
23
24  > # Check the decomposition of SStr
25  > myaov = aov( Y ~ treat )
26  > summary.lm(myaov)
27  Call:
28  aov(formula = Y ~ treat)
29  Residuals:
30      Min     1Q  Median     3Q     Max
31  -1.8500 -0.7125  0.1750  0.4625  1.8500
32  Coefficients:
33              Estimate Std. Error t value Pr(>|t|)
34  (Intercept)   2.2300     0.2379   9.375 1.16e-07 ***
35  treat.L      -1.0356     0.5319  -1.947  0.07050 .
36  treat.Q      -0.8886     0.5319  -1.671  0.11551
37  treat.C      -0.4981     0.5319  -0.936  0.36391
38  treat^4       2.1544     0.5319   4.050  0.00105 **
39  ---
40  Residual standard error: 1.064 on 15 degrees of freedom
41  Multiple R-squared:  0.614,    Adjusted R-squared:  0.5111
42  F-statistic: 5.966 on 4 and 15 DF,  p-value: 0.004442
43
44  > MSE = summary(myaov)[[1]][["Mean Sq"]][2]
45  > t.stat = summary.lm(myaov)[[4]][-1,3]
46  > SStr.for.contrast = t.stat^2 * MSE
47  > SStr.for.contrast
48   treat.L  treat.Q  treat.C  treat^4
49   4.29025  3.15875  0.99225 18.56575
50  > sum(SStr.for.contrast)  ## Note: SStr = 27.007
51  [1] 27.007
```

### Ⓡ: (3) contrasts using `contr.sum()` function

```
1   > contrasts(treat) = contr.sum(5)
2   > contrasts(treat)
3     [,1] [,2] [,3] [,4]
4   A    1    0    0    0
5   B    0    1    0    0
6   C    0    0    1    0
7   D    0    0    0    1
8   E   -1   -1   -1   -1
9   >
10  > # Check the validity of contrasts (sum=0)
11  > cnt = contrasts(treat)
12  > apply(cnt, 2, sum)
13  [1] 0 0 0 0
14
15
16  > # Check the orthogonality (not orthogonal as shown below)
17  > crossprod(cnt)
18       [,1] [,2] [,3] [,4]
19  [1,]    2    1    1    1
20  [2,]    1    2    1    1
21  [3,]    1    1    2    1
22  [4,]    1    1    1    2
23
24  > # Check the decomposition of SStr
25  > myaov = aov( Y ~ treat )
26  > summary.lm(myaov)
27  Call:
28  aov(formula = Y ~ treat)
29  Residuals:
30      Min     1Q  Median     3Q     Max
31  -1.8500 -0.7125  0.1750  0.4625  1.8500
32  Coefficients:
33              Estimate Std. Error t value Pr(>|t|)
34  (Intercept)   2.2300     0.2379   9.375 1.16e-07 ***
35  treat1        0.5950     0.4757   1.251 0.230213
36  treat2       -0.7800     0.4757  -1.640 0.121901
37  treat3        2.0200     0.4757   4.246 0.000704 ***
38  treat4       -0.8050     0.4757  -1.692 0.111293
```

```
39  ---
40  Residual standard error: 1.064 on 15 degrees of freedom
41  Multiple R-squared:  0.614,    Adjusted R-squared:  0.5111
42  F-statistic: 5.966 on 4 and 15 DF,  p-value: 0.004442
43
44  > MSE = summary(myaov)[[1]][["Mean Sq"]][2]
45  > t.stat = summary.lm(myaov)[[4]][-1,3]
46  > SStr.for.contrast = t.stat^2 * MSE
47  > SStr.for.contrast
48     treat1     treat2     treat3     treat4
49   1.770125   3.042000  20.402000   3.240125
50  > sum(SStr.for.contrast)  ## Note: SStr = 27.007
51  [1] 28.45425
```

$\|$

## 6.3  Contrasts – different version

When samples are unbalanced, great care should be taken. Contrast constants and their orthogonality are defined in Definitions 2 and 3. However, some textbooks use a different version as below.

---

**Definition 5.** We denote $\boldsymbol{\theta} = (\mu_1, \mu_2, \ldots, \mu_r)$ be a set of parameters (or statistics) and $\mathbf{c} = (c_1, c_2, \ldots, c_r)$ be a collection of known constants. We assume that $\mu_i$ is estimated with a sample of size $n_i$. Then $\mathbf{c} = (c_1, c_2, \ldots, c_r)$ is called *contrast constants* and its linear combination below is called a *contrast*

$$\Gamma = \sum_{i=1}^{r} n_i c_i \mu_i,$$

where $\sum_{i=1}^{r} n_i c_i = 0$.

---

**Remark 8.**  We can notice the following.

1. The condition $\sum_{i=1}^{r} n_i c_i = 0$ is similar to the condition in (3).

2. The textbook uses $\sum_{i=1}^{r} c_i = 0$ (Definition 2) condition in 정의 3.2 on Page 68.

3. It is easily seen that $a_i$ in Definition 2 and $c_i$ in Definition 5 has the relation $a_i = n_i c_i$.

$\triangle$

**Definition 6.** Let $\overline{Y}_{i\bullet}$ be the estimator of the $i$ treatment $\mu_i$ based on $n_i$ observations and $\mathbf{c} = (c_1, c_2, \ldots, c_r)$ and $\mathbf{d} = (d_1, d_2, \ldots, d_r)$ be contrast constants satisfying

$$\sum_{i=1}^{r} n_i c_i d_i = 0.$$

Then two contrasts below are called *orthogonal*

$$\Gamma_{\mathbf{a}} = \sum_{i=1}^{r} n_i c_i \mu_i \ \text{ and } \ \Gamma_{\mathbf{b}} = \sum_{i=1}^{r} n_i d_i \mu_i.$$

**Remark 9.**   It should be noted that the textbook uses $\sum_{i=1}^{r} n_i c_i d_i = 0$ (Definition 6) in 정의 3.3 on Page 70. △

**Lemma 6.**   *Under the assumption that $Y_{ij} = \mu_i + \epsilon_{ij}$ with $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$, we have*

$$\mathrm{Cov}\left( \sum_{i=1}^{r} n_i c_i \overline{Y}_{i\bullet}, \sum_{j=1}^{r} n_i d_j \overline{Y}_{j\bullet} \right) = \sigma^2 \sum_{i=1}^{r} n_i c_i d_i.$$

*Proof.*   The proof is very similar to that of Lemma 5.                                              □

Analogous to (29), we have $n_i \overline{Y}_{i\bullet} \sim N\left( n_i \mu_i, n_i \sigma^2 \right)$. Thus, for any contrast constants $\mathbf{c} = (c_1, c_2, \ldots, c_r)$, $\sum_{i=1}^{r} n_i c_i \overline{Y}_{i\bullet}$ is also normally distributed with

$$E\left( \sum_{i=1}^{r} n_i c_i \overline{Y}_{i\bullet} \right) = \sum_{i=1}^{r} n_i c_i \mu_i \ \text{ and } \ \mathrm{Var}\left( \sum_{i=1}^{r} n_i c_i \overline{Y}_{i\bullet} \right) = \sigma^2 \sum_{i=1}^{r} n_i c_i^2.$$

Standardizing $\sum_{i=1}^{r} n_i c_i \overline{Y}_{i\bullet}$, we have

$$\frac{\sum_{i=1}^{r} n_i c_i \overline{Y}_{i\bullet} - \sum_{i=1}^{r} n_i c_i \mu_i}{\sqrt{\sigma^2 \sum_{i=1}^{r} n_i c_i^2}} \sim N(0, 1). \tag{35}$$

Note that (35) defines a pivot but it includes a nuisance parameter $\sigma^2$. As we Studentized (30), we can also do (35) so that we have

$$\frac{\sum_{i=1}^{r} n_i c_i \overline{Y}_{i\bullet} - \sum_{i=1}^{r} n_i c_i \mu_i}{\sqrt{\mathrm{MSE} \sum_{i=1}^{r} n_i c_i^2}} \sim t(\mathrm{df} = N - r).$$

As we decompose the treatment sum of squares (SStr) into $r - 1$ independent random variables with all one degree of freedom with orthogonal (uncorrelated) contrasts in Definition 2, we can also do this decomposition with orthogonal (uncorrelated) contrasts in

Definition 5. Denote these sets by $\mathbf{c}^{(\ell)} = (c_1^{(\ell)}, c_2^{(\ell)}, \ldots, c_r^{(\ell)})$ for $\ell = 1, 2, \ldots, r-1$ which satisfy Definition 6, that is,

$$\sum_{i=1}^{r} n_i c_i^{(\ell)} c_i^{(\ell')} = 0,$$

for all $\ell \neq \ell'$. Analogous to (33), we can decompose SStr in (15) into contrast sums of squares as follows

$$\sum_{i=1}^{r} n_i (\overline{Y}_{i\bullet} - \overline{Y}_{\bullet\bullet})^2 = \sum_{\ell=1}^{r-1} \frac{\left( \sum_{i=1}^{r} n_i c_i^{(\ell)} \overline{Y}_{i\bullet} \right)^2}{\sum_{i=1}^{r} n_i \left( c_i^{(\ell)} \right)^2}.$$

**Example 6.** Revisit Example 2. In this example, we calculated $\text{SStr} = 22.997 \approx 23$.

Ⓡ: Reading Data

```
1  > Y1 = c(25.7, 27.2, 29.9, 28.5, 29.4, 28.3)          # Blue
2  > Y2 = c(26.8, 27.9, 23.7, 25, 26.3, 24.8, 25.7, 24.5) # Brown
3  > Y3 = c(26.4, 24.2, 28.0, 26.9, 29.1)                # Green
4  > Y = c(Y1, Y2, Y3)
5  > n1 = length(Y1); n2 = length(Y2); n3 = length(Y3)
6  > ni = c(n1,n2,n3)
```

Ⓡ: SStr, MSE, etc.

```
1  > # ni and Yi.bar
2  > Y1bar = mean(Y1); Y2bar = mean(Y2); Y3bar = mean(Y3)
3  > Yi.bar = c(Y1bar, Y2bar, Y3bar)
4  > rbind(ni, Yi.bar)
5             [,1]     [,2]   [,3]
6  ni       6.00000  8.0000  5.00
7  Yi.bar  28.16667 25.5875 26.92
8  >
9  > # SStr
10 > SStr = sum( ni*(Yi.bar-mean(Y))^2 )
11 > SStr
12 [1] 22.99729
13 >
14 > # SSE
15 > df = n1+n2+n3-3
16 > SSE = sum((Y1-Y1bar)^2) + sum((Y2-Y2bar)^2) + sum((Y3-Y3bar)^2)
17 > MSE = SSE / df
18 > MSE
19 [1] 2.39438
```

Ⓡ: (1) Contrasts $a_i$ (correct)

```
1  > a1 = c(-0.6, 1.1, -0.5);  a2 = c(-1, 0, 1)
2  > sum(a1)           # zero
3  [1] 1.110223e-16
4  > sum(a2)           # zero
5  [1] 0
6  > sum( a1*a2 / ni ) # zero
7  [1] -1.387779e-17
8  > sum( ni*a1*a2 )   # not zero (does not satisfy the textbook definition)
9  [1] 1.1
10 >
11 > SSa1 = sum( a1*Yi.bar )^2 / sum( a1^2/ni )
12 > SSa2 = sum( a2*Yi.bar )^2 / sum( a2^2/ni )
13 > c(SSa1,SSa2, SSa1+SSa2) # Correct decomposition (SStr=22.997)
14 [1] 18.758618  4.238667 22.997285
15 >
```

```
16  > # Double check with t-test statistics
17  > T1 = sum(a1*Yi.bar) / sqrt( sum(MSE*a1^2*(1/ni)) )
18  > T2 = sum(a2*Yi.bar) / sqrt( sum(MSE*a2^2*(1/ni)) )
19  > c(T1^2*MSE, T2^2*MSE, T1^2*MSE+T2^2*MSE)     ## SStr=22.997
20  [1] 18.758618  4.238667 22.997285
21  > # Correct decomposition: 18.758618 + 4.238667 = 22.997285 (SStr)
```

### Ⓡ: (2) Contrasts $c_i$ from the textbook (wrong)

```
1   > c1 = c(-0.5, 1.1, -0.6);  c2 = c(-1, 0, 1)
2   > sum(c1)            # zero (satisfies Definition 3.2 in the textbook)
3   [1] 1.110223e-16
4   > sum(c2)            # zero (satisfies Definition 3.2 in the textbook)
5   [1] 0
6   > sum( ni*c1*c2 )    # zero (satisfies Definition 3.3 in the textbook)
7   [1] 0
8   >
9   > SSc1 = sum( c1*Yi.bar )^2 / sum( c1^2/ni )
10  > SSc2 = sum( c2*Yi.bar )^2 / sum( c2^2/ni )
11  > c(SSc1,SSc2, SSc1+SSc2) # Wrong decomposition (SStr=22.997)
12  [1] 16.474121  4.238667 20.712787
13  >
14  > # Double check with t-test statistics
15  > T1 = sum(c1*Yi.bar) / sqrt( sum(MSE*c1^2*(1/ni)) )
16  > T2 = sum(c2*Yi.bar) / sqrt( sum(MSE*c2^2*(1/ni)) )
17  > c(T1^2*MSE, T2^2*MSE, T1^2*MSE+T2^2*MSE)
18  [1] 16.474121  4.238667 20.712787
19  > # Wrong decomposition: 16.474121 + 4.238667 = 20.712787  (not 22.997)
```

### Ⓡ: (3) Contrasts $c_i$ from the handout (correct)

```
1   > a1 = c(-0.6, 1.1, -0.5);  a2 = c(-1.0, 0.0,  1.0)
2   > c1 = a1/ni;               c2 = a2/ni
3   > sum(ni*c1)          # zero (satisfies the definition in the handout)
4   [1] 1.110223e-16
5   > sum(ni*c2)          # zero (satisfies the definition in the handout)
6   [1] 0
7   > sum(ni*c1*c2)       # zero (satisfies the definition in the handout)
8   [1] -1.387779e-17
9   >
10  > SSc1 = sum( c1*ni*Yi.bar )^2 / sum( ni*c1^2 )
11  > SSc2 = sum( c2*ni*Yi.bar )^2 / sum( ni*c2^2 )
12  > c(SSc1, SSc2, SSc1+SSc2)
13  [1] 18.758618  4.238667 22.997285
14  > # Correct decomposition: 18.758618 + 4.238667 = 22.997285 (SStr)
```

∥

**Example 7.**   We solve Example 6 again using the `aov()` function. As shown in the R programs below, the orthogonal contrasts ($a_i$) in Definition 3 fail to decompose the SStr. Only the orthogonal contrasts ($c_i$) satisfying Definitions 5 and 6 work properly when the `aov()` function is used. Note that both $a_i$ or $c_i$ decomposed the SStr successfully in Example 6. However, the contrasts $c_i$ based on 정의 3.2 of the textbook fail to decompose the SStr in both Examples 6 and 7.

### Ⓡ: Reading Data

```
1   > Y1 = c(25.7, 27.2, 29.9, 28.5, 29.4, 28.3)           # Blue
2   > Y2 = c(26.8, 27.9, 23.7, 25, 26.3, 24.8, 25.7, 24.5) # Brown
3   > Y3 = c(26.4, 24.2, 28.0, 26.9, 29.1)                 # Green
4   > Y = c(Y1, Y2, Y3)
```

```
5  > n1 = length(Y1); n2 = length(Y2); n3 = length(Y3)
6  > ni = c(n1,n2,n3)
7  > color = factor( rep(c("Blue","Brown","Green"), ni) )
```

## Ⓡ: (1) Contrasts $a_i$ (wrong)

```
1   > # Contrasts
2   > a1 = c(-0.6, 1.1, -0.5);  a2 = c(-1, 0, 1)
3   > sum(a1)            # zero
4   [1] 1.110223e-16
5   > sum(a2)            # zero
6   [1] 0
7   > sum( a1*a2 / ni ) # zero
8   [1] -1.387779e-17
9
10  > contrasts(color) = cbind(a1,a2)
11  > myaov = aov( Y ~ color )
12  > myaov
13  Call:
14     aov(formula = Y ~ color)
15  Terms:
16                     color Residuals
17  Sum of Squares  22.99729  38.31008
18  Deg. of Freedom        2        16
19
20  Residual standard error: 1.547378
21  Estimated effects may be unbalanced
22
23  > summary(myaov)
24             Df Sum Sq Mean Sq F value Pr(>F)
25  color       2  23.00  11.499   4.802 0.0232 *
26  Residuals  16  38.31   2.394
27  ---
28
29  > # t-test statistics
30  > summary.lm(myaov)
31  Call:
32  aov(formula = Y ~ color)
33  Residuals:
34      Min      1Q  Median      3Q     Max
35  -2.7200 -0.8771  0.1125  1.1462  2.3125
36
37  Coefficients:
38             Estimate Std. Error t value Pr(>|t|)
39  (Intercept) 26.8914     0.3617  74.354   <2e-16 ***
40  colora1     -1.1854     0.4365  -2.715   0.0153 *
41  colora2     -0.5641     0.4703  -1.199   0.2478
42  ---
43  Residual standard error: 1.547 on 16 degrees of freedom
44  Multiple R-squared:  0.3751,    Adjusted R-squared:  0.297
45  F-statistic: 4.802 on 2 and 16 DF,  p-value: 0.02325
46
47  > # Decomposition
48  > MSE = summary(myaov)[[1]][["Mean Sq"]][2]
49  > t.stat = summary.lm(myaov)[[4]][-1,3]
50  > SStr.for.contrast = t.stat^2 * MSE
51  > SStr.for.contrast
52    colora1    colora2
53  17.655157   3.444492
54  > sum(SStr.for.contrast)
55  [1] 21.09965
56  > # Wrong decomposition: 17.655157  3.444492 = 21.09965 (not 22.997)
```

## Ⓡ: (2) Contrasts $c_i$ from the textbook (wrong)

```
1  > # Contrasts
2  > c1 = c(-0.5, 1.1, -0.6);  c2 = c(-1, 0, 1)
3  > sum(c1)            # zero (satisfies Definition 3.2 in the textbook)
4  [1] 1.110223e-16
5  > sum(c2)            # zero (satisfies Definition 3.2 in the textbook)
6  [1] 0
```

```
 7  > sum( ni*c1*c2 )    # zero (satisfies Definition 3.3 in the textbook)
 8  [1] 0
 9
10  > contrasts(color) = cbind(c1,c2)
11  > myaov = aov( Y ~ color )
12  > myaov
13  Call:
14      aov(formula = Y ~ color)
15  Terms:
16                        color Residuals
17  Sum of Squares  22.99729   38.31008
18  Deg. of Freedom        2         16
19
20  Residual standard error: 1.547378
21  Estimated effects may be unbalanced
22  > summary(myaov)
23             Df Sum Sq Mean Sq F value Pr(>F)
24  color       2  23.00   11.499   4.802 0.0232 *
25  Residuals  16  38.31    2.394
26  ---
27
28  > # t-test statistics
29  > summary.lm(myaov)
30  Call:
31  aov(formula = Y ~ color)
32  Residuals:
33      Min      1Q  Median      3Q     Max
34  -2.7200 -0.8771  0.1125  1.1462  2.3125
35  Coefficients:
36              Estimate Std. Error t value Pr(>|t|)
37  (Intercept)  26.8914     0.3617  74.354   <2e-16 ***
38  colorc1      -1.1854     0.4365  -2.715   0.0153 *
39  colorc2      -0.6826     0.4677  -1.459   0.1638
40  ---
41  Residual standard error: 1.547 on 16 degrees of freedom
42  Multiple R-squared:  0.3751,     Adjusted R-squared:  0.297
43  F-statistic: 4.802 on 2 and 16 DF,  p-value: 0.02325
44
45  > # Decomposition
46  > MSE = summary(myaov)[[1]][["Mean Sq"]][2]
47  > t.stat = summary.lm(myaov)[[4]][-1,3]
48  > SStr.for.contrast = t.stat^2 * MSE
49  > SStr.for.contrast
50    colorc1    colorc2
51  17.655157   5.100057
52  > sum(SStr.for.contrast)
53  [1] 22.75521
54  > # Wrong decomposition: 17.655157 + 5.100057 = 22.75521 (not 22.997)
```

## Ⓡ: (3) Contrasts $c_i$ from the handout (correct)

```
 1  > # Try the contrasts below
 2  > a1 = c(-0.6, 1.1, -0.5);  a2 = c(-1.0, 0.0,  1.0)
 3  > c1 = a1/ni;               c2 = a2/ni
 4  > sum(ni*c1)         # zero (satisfies the definition in the handout)
 5  [1] 1.110223e-16
 6  > sum(ni*c2)         # zero (satisfies the definition in the handout)
 7  [1] 0
 8  > sum(ni*c1*c2)      # zero (satisfies the definition in the handout)
 9  [1] -1.387779e-17
10
11  > contrasts(color) = cbind(c1,c2)
12  > myaov = aov( Y ~ color )
13  > myaov
14  Call:
15      aov(formula = Y ~ color)
16  Terms:
17                        color Residuals
18  Sum of Squares  22.99729   38.31008
19  Deg. of Freedom        2         16
20
```

```
21   Residual standard error: 1.547378
22   Estimated effects are balanced
23   > summary(myaov)
24               Df Sum Sq Mean Sq F value Pr(>F)
25   color        2  23.00  11.499   4.802 0.0232 *
26   Residuals   16  38.31   2.394
27   ---
28
29   > # t-test statistics
30   > summary.lm(myaov)
31   Call:
32   aov(formula = Y ~ color)
33   Residuals:
34       Min      1Q  Median      3Q     Max
35   -2.7200 -0.8771  0.1125  1.1462  2.3125
36
37   Coefficients:
38               Estimate Std. Error t value Pr(>|t|)
39   (Intercept)   26.753      0.355  75.361   <2e-16 ***
40   colorc1       -8.474      3.027  -2.799   0.0129 *
41   colorc2       -3.400      2.555  -1.331   0.2020
42   ---
43   Residual standard error: 1.547 on 16 degrees of freedom
44   Multiple R-squared:  0.3751,     Adjusted R-squared:  0.297
45   F-statistic: 4.802 on 2 and 16 DF,  p-value: 0.02325
46
47   > # Decomposition
48   > MSE = summary(myaov)[[1]][["Mean Sq"]][2]
49   > t.stat = summary.lm(myaov)[[4]][-1,3]
50   > SStr.for.contrast = t.stat^2 * MSE
51   > SStr.for.contrast
52     colorc1   colorc2
53   18.758618  4.238667
54   > sum(SStr.for.contrast)
55   [1] 22.99729
56   > # Correct decomposition: 18.758618 + 4.238667 = 22.997285 (SStr)
```

$\parallel$

## 6.4  Comparing pairs of treatment means

Using contrasts, we compare any differences in treatment means. However, in many practice, it is very difficult to determine which contrasts are to be used. One appealing idea proposed by R. A. Fisher is comparing only *pairs* of treatment means which can be done with the contrasts $\Gamma = \mu_\ell - \mu_m$ for $\ell \neq m$. The Fisher method compares all the pairs of treatment means with the significance level $\alpha$ for each individual pair. The underlying theory on this method was already investigated in Section 1.2. Under $H_0 : \mu_\ell = \mu_m$, the $t$-test statistic in (13) becomes

$$\frac{\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet}}{\sqrt{\text{MSE}\left(\frac{1}{n_\ell} + \frac{1}{n_m}\right)}} \sim t(\text{df} = N - r).$$

Also, we obtained the endpoints for the $100(1-\alpha)\%$ confidence interval of $\mu_\ell - \mu_m$ in (14). Notice that the deviation from the center of the interval is given by

$$t(1 - \tfrac{\alpha}{2}; N - r)\sqrt{\text{MSE}\left(\frac{1}{n_\ell} + \frac{1}{n_m}\right)}.$$

which is called the least significance difference (LSD) and we denote it by $\text{LSD}_{\ell m}$. To perform this method, we simply compare the observed difference between each pair of means to the corresponding LSD. That is, if $|\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet}| > \text{LSD}_{\ell m}$, then we can conclude that the treatment means $\mu_\ell$ and $\mu_m$ are different.

These days, with the advent of powerful and accessible computers, it is quite easy to calculate $p$-value. Thus, we also carry out this method by comparing the $p$-values of the $t$-test statistic for each pair of means. In general, $p$-values are preferred because they provide better information than test statistics. Let $T_{\ell m} = (\overline{Y}_{\ell\bullet} - \overline{Y}_{m\bullet})/\sqrt{\text{MSE}(1/n_\ell + 1/n_m)}$. Then the corresponding $p$-value is calculated as

$$2 \cdot \left\{ 1 - F_{N-r}\big(|T_{\ell m}|\big) \right\},$$

where $F_\nu(\,\cdot\,)$ is the cumulative distribution function of the $t$-distribution with $\nu$ degrees of freedom. Note that $1 - F_\nu(\,\cdot\,)$ is easily calculated with `pt(q, df=`$\nu$`, lower.tail=FALSE)` function in R language.

**Remark 10.**    There are several other methods for this pairwise comparison. A natural question is: *Which method is recommended?* There is no clear answer for this question. Carmer and Swanson (1973) have carried out extensive Monte Carlo simulations whose results show that the Fisher's least significant difference (LSD) method is very effective. But, the Fisher's LSD method can not control the overall error rate (say, at the selected level $\alpha$) which is also known as *family* or *experiment-wise* error rate. Thus, the overall error rate is needed, refer to other methods such as Scheffé (1953) and Tukey (1953). $\triangle$

**Example 8.**    We solve Example 6 again to use the Fisher LSD procedure.

Ⓡ: Reading Data

```
1  > Y1 = c(25.7, 27.2, 29.9, 28.5, 29.4, 28.3)        # Blue
2  > Y2 = c(26.8, 27.9, 23.7, 25, 26.3, 24.8, 25.7, 24.5) # Brown
3  > Y3 = c(26.4, 24.2, 28.0, 26.9, 29.1)               # Green
4  > Y = c(Y1, Y2, Y3)
5  > n1 = length(Y1); n2 = length(Y2); n3 = length(Y3)
6  > df = n1+n2+n3-3
```

Ⓡ: Using $t$-test statistics

```
1  > # Don't use: t.test(Y1,Y2, var.equal=TRUE)
2  > # MSE
3  > SSE = sum((Y1-mean(Y1))^2) + sum((Y2-mean(Y2))^2) + sum((Y3-mean(Y3))^2)
4  > MSE = SSE / df
5  > MSE
6  [1] 2.39438
7  >
8  > T12 = (mean(Y1)-mean(Y2))/sqrt(MSE*(1/n1+1/n2))
9  > p12 = 2*pt(abs(T12), df=df, lower.tail=FALSE)
10 > c(T12, p12) # Blue vs. Brown
11 [1] 3.086309326 0.007079982
12 >
13 > T13 = (mean(Y1)-mean(Y3))/sqrt(MSE*(1/n1+1/n3))
14 > p13 = 2*pt(abs(T13), df=df, lower.tail=FALSE)
15 > c(T13, p13) # Blue vs. Green
16 [1] 1.3305098 0.2020033
17 >
18 > T23 = (mean(Y2)-mean(Y3))/sqrt(MSE*(1/n2+1/n3))
19 > p23 = 2*pt(abs(T23), df=df, lower.tail=FALSE)
20 > c(T23, p23) # Brown vs. Green
21 [1] -1.5105287  0.1504046
22 >
23 > # Compare
24 > # t.test(Y1,Y2, var.equal=TRUE)$p.value  vs. p12
25 > # t.test(Y1,Y3, var.equal=TRUE)$p.value  vs. p13
26 > # t.test(Y2,Y3, var.equal=TRUE)$p.value  vs. p23
```

Ⓡ: Using `pairwise.t.test()` function

```
1  > color = factor( rep(c("Blue","Brown","Green"), c(n1,n2,n3)) )
2  > pairwise.t.test(Y,color, p.adjust="none", pool.sd=TRUE)
3
4          Pairwise comparisons using t tests with pooled SD
5
6  data:  Y and color
7
8         Blue   Brown
9  Brown 0.0071 -
10 Green 0.2020 0.1504
11
12 P value adjustment method: none
```

$\|$

**Example 9.**    We consider Example 4 again to use the Fisher LSD procedure.

Ⓡ: Reading Data

```
1  > Y1 = c(2.4, 2.7, 3.1, 3.1)
2  > Y2 = c(0.7, 1.6, 1.7, 1.8)
3  > Y3 = c(2.4, 3.1, 5.4, 6.1)
4  > Y4 = c(0.3, 0.3, 2.4, 2.7)
5  > Y5 = c(0.5, 0.9, 1.4, 2.0)
6  > Y = c(Y1, Y2, Y3, Y4, Y5)
7  > n1=length(Y1);n2=length(Y2);n3=length(Y3);n4=length(Y4);n5=length(Y5)
8  > df = n1+n2+n3+n4+n5-5
```

Ⓡ: Using $t$-test statistics

```
1  > SSE = sum((Y1-mean(Y1))^2) + sum((Y2-mean(Y2))^2) +
      sum((Y3-mean(Y3))^2) + sum((Y4-mean(Y4))^2) + sum((Y5-mean(Y5))^2)
2  > MSE = SSE / df
3  > MSE
4  [1] 1.131667
5
```

```
6    > # 1 vs. {2,3,4,5}
7    > T12 = (mean(Y1)-mean(Y2))/sqrt(MSE*(1/n1+1/n2))
8    > p12 = 2*pt(abs(T12), df=df, lower.tail=FALSE)
9    > c(T12, p12)
10   [1] 1.82792526 0.08751812
11
12   > T13 = (mean(Y1)-mean(Y3))/sqrt(MSE*(1/n1+1/n3))
13   > p13 = 2*pt(abs(T13), df=df, lower.tail=FALSE)
14   > c(T13, p13)
15   [1] -1.89439527  0.07761809
16
17   > T14 = (mean(Y1)-mean(Y4))/sqrt(MSE*(1/n1+1/n4))
18   > p14 = 2*pt(abs(T14), df=df, lower.tail=FALSE)
19   > c(T14, p14)
20   [1] 1.86116026 0.08243548
21
22   > T15 = (mean(Y1)-mean(Y5))/sqrt(MSE*(1/n1+1/n5))
23   > p15 = 2*pt(abs(T15), df=df, lower.tail=FALSE)
24   > c(T15, p15)
25   [1] 2.16027531 0.04734312
26
27   > # 2 vs. {3,4,5}
28   > T23 = (mean(Y2)-mean(Y3))/sqrt(MSE*(1/n2+1/n3))
29   > p23 = 2*pt(abs(T23), df=df, lower.tail=FALSE)
30   > c(T23, p23)
31   [1] -3.722320527  0.002043516
32   >
33   > T24 = (mean(Y2)-mean(Y4))/sqrt(MSE*(1/n2+1/n4))
34   > p24 = 2*pt(abs(T24), df=df, lower.tail=FALSE)
35   > c(T24, p24)
36   [1] 0.0332350 0.9739254
37   >
38   > T25 = (mean(Y2)-mean(Y5))/sqrt(MSE*(1/n2+1/n5))
39   > p25 = 2*pt(abs(T25), df=df, lower.tail=FALSE)
40   > c(T25, p25)
41   [1] 0.332350 0.744224
42
43   > # 3 vs. {4,5}
44   > T34 = (mean(Y3)-mean(Y4))/sqrt(MSE*(1/n3+1/n4))
45   > p34 = 2*pt(abs(T34), df=df, lower.tail=FALSE)
46   > c(T34, p34)
47   [1] 3.755555531 0.001909124
48   >
49   > T35 = (mean(Y3)-mean(Y5))/sqrt(MSE*(1/n3+1/n5))
50   > p35 = 2*pt(abs(T35), df=df, lower.tail=FALSE)
51   > c(T35, p35)
52   [1] 4.054670574 0.001037412
53
54   > # 4 vs. 5
55   > T45 = (mean(Y4)-mean(Y5))/sqrt(MSE*(1/n4+1/n5))
56   > p45 = 2*pt(abs(T45), df=df, lower.tail=FALSE)
57   > c(T45, p45)
58   [1] 0.2991150 0.7689566
59
60   > cbind( c(p12,p13,p14,p15), c(NA,p23,p24,p25), c(NA,NA,p34,p35),
         c(NA,NA,NA,p45))
61             [,1]        [,2]        [,3]      [,4]
62   [1,] 0.08751812         NA          NA        NA
63   [2,] 0.07761809 0.002043516         NA        NA
64   [3,] 0.08243548 0.973925402 0.001909124        NA
65   [4,] 0.04734312 0.744223954 0.001037412 0.7689566
```

### Ⓡ: Using `pairwise.t.test()` function

```
1    > treat = factor( rep(1:5, c(n1,n2,n3,n4,n5)) )
2    > pairwise.t.test(Y, treat, p.adjust="none", pool.sd=TRUE)
3
4            Pairwise comparisons using t tests with pooled SD
5    data:  Y and treat
6      1      2      3      4
7    2 0.0875 -      -      -
```

```
 8  3 0.0776 0.0020 -      -
 9  4 0.0824 0.9739 0.0019 -
10  5 0.0473 0.7442 0.0010 0.7690
```

Where there are many groups, it is more convenient to sort the factors based on their corresponding treatment means from the smallest to the largest. Then, it is easier to cluster treatments.

Ⓡ: Using `pairwise.t.test()` function after soring the group means

```
 1 > treat = factor( rep(1:5, c(n1,n2,n3,n4,n5)) )
 2
 3 > by(Y, treat, mean)   # sample means for each treatment
 4 treat: 1
 5 [1] 2.825
 6 ----------------------------------------------------------
 7 treat: 2
 8 [1] 1.45
 9 ----------------------------------------------------------
10 treat: 3
11 [1] 4.25
12 ----------------------------------------------------------
13 treat: 4
14 [1] 1.425
15 ----------------------------------------------------------
16 treat: 5
17 [1] 1.2
18 > gr.order = order(by(Y, treat, mean)) # the order of the means
19 > gr.order
20 [1] 5 4 2 1 3
21
22 > treat2 = factor(treat, levels=gr.order)
23 > levels(treat2)
24 [1] "5" "4" "2" "1" "3"
25 > pairwise.t.test(Y, treat2, p.adjust="none", pool.sd=TRUE)
26
27         Pairwise comparisons using t tests with pooled SD
28 data:  Y and treat2
29   5      4      2      1
30 4 0.7690 -      -      -
31 2 0.7442 0.9739 -      -
32 1 0.0473 0.0824 0.0875 -
33 3 0.0010 0.0019 0.0020 0.0776
```

$\|$

# 7   Note on heteroscedasticity case

For testing the effect of the treatments, we assumed the model $Y_{ij} = \mu_i + \epsilon_{ij}$ where $\epsilon_{ij} \overset{iid}{\sim} N(0, \sigma^2)$. This model assumes that the variances of all the samples are equal, which is called *homoscedasticity*. If this homoscedasticity is not satisfied so that $\epsilon_{ij} \sim N(0, \sigma_i^2)$, we have to use other methods. In the statistics literature, several methods are recommended. However, there is no clear winner for this heteroscedasticity case. Some well-known methods are (i) Box-Cox transform, (ii) Welch-type ANOVA, (iii) weighted linear regression model, (iv) bootstrap, etc.

# References

Carmer, S. G. and Swanson, M. R. (1973). An evaluation of ten pairwise multiple comparison procedures by monte carlo methods. *Journal of the American Statistical Association*, 68(341):66–74.

Casella, G. and Berger, R. L. (2002). *Statistical Inference*. Duxbury, Pacific Grove, CA, second edition.

Cochran, W. G. (1934). The distribution of quadratic forms in a normal system with applications to the analysis of variance. *Proceedings of the Cambridge Philosophical Society*, 30:178–191.

Hogg, R. V., Tanis, E. A., and Zimmerman, D. L. (2015). *Probability and Statistical Inference*. Pearson, 9th edition.

Kendall, M. and Stuart, A. (1979). *The Advanced Theory of Statistics*, volume 2. Charles Griffin, fourth edition.

Ross, S. M. (2014). *A First Course in Probability*. Prentice Hall, 9th edition.

Scheffé, H. (1953). A method for judging all contrasts in the analysis of variance. *Biometrika*, 40:87–104.

Smith, J. M. and Misiak, H. (1973). The effect of iris color on critical flicker frequency. *The Journal of General Psychology*, 89:91–95.

Tukey, J. W. (1953). The problem of multiple comparisons. Unpublished notes, Princeton University.