

Math 362: Mathematical Statistics II

Le Chen

le.chen@emory.edu

Emory University
Atlanta, GA

Last updated on April 24, 2021

2021 Spring

Chapter 12. The Analysis of Variance

§ 12.1 Introduction

§ 12.2 The F Test

§ 12.3 Multiple Comparisons: Turkey's Method

§ 12.4 Testing Subhypotheses with Contrasts

Chapter 12. The Analysis of Variance

§ 12.1 Introduction

§ 12.2 The F Test

§ 12.3 Multiple Comparisons: Turkey's Method

§ 12.4 Testing Subhypotheses with Contrasts

Model assumptions

1. Independence of observations
2. Normality
3. Homogeneity of variances



Assume:

$\forall j = 1, \dots, k, \forall i = 1, \dots, n_j,$

1. Y_{ij} are independent.
2. $Y_{ij} \sim N(\mu_j, \sigma^2)$



Assume:

$\forall j = 1, \dots, k, \forall i = 1, \dots, n_j,$

$$Y_{ij} = \mu_j + \epsilon_{ij}$$

1. ϵ_{ij} are independent.
2. $\epsilon_{ij} \sim N(0, \sigma^2)$

Table 12.1.1

	<i>Treatment Level</i>			
	1	2	...	k
	Y_{11}	Y_{12}		Y_{1k}
	Y_{21}	Y_{22}		
	\vdots	\vdots	...	\vdots
	$Y_{n_1 1}$	$Y_{n_2 2}$		$Y_{n_k k}$
Sample sizes:	n_1	n_2	...	n_k
Sample totals:	$T_{\cdot 1}$	$T_{\cdot 2}$		$T_{\cdot k}$
Sample means:	$\bar{Y}_{\cdot 1}$	$\bar{Y}_{\cdot 2}$		$\bar{Y}_{\cdot k}$
True means:	μ_1	μ_2		μ_k

Likelihood ratio test

1. The parameter spaces are

$$\Omega = \{(\mu_1, \dots, \mu_k, \sigma^2) : -\infty < \mu_1, \dots, \mu_k < \infty, \sigma^2 > 0\}$$

$$\omega = \{(\mu_1, \dots, \mu_k, \sigma^2) : -\infty < \mu_1 = \dots = \mu_k < \infty, \sigma^2 > 0\}$$

2. The likelihood functions are

$$L(\omega) = \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \mu)^2 \right\}$$

$$L(\Omega) = \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \mu_j)^2 \right\}$$

3. Now

$$\frac{\partial \ln L(\omega)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{j=1}^k \sum_{i=1}^{\eta_j} (y_{ij} - \mu)$$

$$\frac{\partial \ln L(\omega)}{\partial (\sigma^2)} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{j=1}^k \sum_{i=1}^{\eta_j} (y_{ij} - \mu)^2$$

Setting the above derivatives to zero, the solutions for μ and σ^2 are,

$$\frac{1}{n} \sum_{j=1}^k \sum_{i=1}^{\eta_j} y_{ij} = \bar{y}_{..}$$

$$\frac{1}{n} \sum_{j=1}^k \sum_{i=1}^{\eta_j} (y_{ij} - \bar{y}_{..})^2 = v$$

3' Similarly,

$$\frac{\partial \ln L(\Omega)}{\partial \mu_j} = \frac{1}{\sigma^2} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \mu_j), \quad j = 1, \dots, k$$

$$\frac{\partial \ln L(\Omega)}{\partial (\sigma^2)} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \mu_j)^2$$

Setting the above derivatives to zero, the solutions for μ_j and σ^2 are,

$$\frac{1}{n_j} \sum_{i=1}^{n_j} y_{ij} = \bar{y}_{\cdot j}$$
$$\frac{1}{n} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j})^2 = w$$

4. Hence,

$$L(\hat{\omega}) = \left(\frac{n}{2\pi \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{..})^2} \right)^{n/2} \exp \left\{ -\frac{n \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{..})^2}{2 \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{..})^2} \right\}$$

$$\parallel$$

$$\left(\frac{n}{2\pi \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{..})^2} \right)^{n/2} e^{-n/2}$$

Similarly,

$$L(\hat{\Omega}) = \left(\frac{n}{2\pi \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{.j})^2} \right)^{n/2} \exp \left\{ -\frac{n \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{.j})^2}{2 \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{.j})^2} \right\}$$

$$\parallel$$

$$\left(\frac{n}{2\pi \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{.j})^2} \right)^{n/2} e^{-n/2}$$

5. Finally,

$$\lambda = \frac{L(\hat{\omega})}{L(\hat{\Omega})} = \left(\frac{\sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j})^2}{\sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot \cdot})^2} \right)^{n/2}$$

⇒ Test statistic:

$$\Lambda = \frac{L(\hat{\omega})}{L(\hat{\Omega})} = \left(\frac{\sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{\cdot j})^2}{\sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{\cdot \cdot})^2} \right)^{n/2}$$

$$\begin{aligned}
SSTOT &:= \sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{..})^2 \\
&= \sum_{j=1}^k \sum_{i=1}^{n_j} \left[(Y_{ij} - \bar{Y}_{.j}) + (\bar{Y}_{.j} - \bar{Y}_{..}) \right]^2 \\
&= \sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{.j})^2 + \text{zero cross term} + \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{Y}_{.j} - \bar{Y}_{..})^2 \\
&= \underbrace{\sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{.j})^2}_{SSE} + \underbrace{\sum_{j=1}^k n_j (\bar{Y}_{.j} - \bar{Y}_{..})^2}_{SSTR}
\end{aligned}$$

⇓

$$\Lambda = \left(\frac{SSE}{SSTOT} \right)^{n/2} = \left(\frac{SSE}{SSE + SSTR} \right)^{n/2} = \left(\frac{1}{1 + SSTR/SSE} \right)^{n/2}$$

6. Critical regions: for some $\lambda_* \in (0, 1)$ close to 0,

$$\begin{aligned}\alpha &= \mathbb{P}(\Lambda \leq \lambda_*) \\ &= \mathbb{P}\left(\frac{1}{1 + SSTR/SSE} \leq \lambda_*^{2/n}\right) \\ &= \mathbb{P}\left(\frac{SSTR}{SSE} \leq \lambda_*^{-2/n} - 1\right) \\ &= \mathbb{P}\left(\frac{SSTR/(k-1)}{SSE/(n-k)} \leq \left(\lambda_*^{-2/n} - 1\right) \frac{n-k}{k-1}\right)\end{aligned}$$

7. We will prove that under H_0 , $\frac{SSTR/(k-1)}{SSE/(n-k)} \sim F$ -distr.
 $df_1 = k - 1, df_2 = n - k$

$$\Rightarrow \left(\lambda_*^{-2/n} - 1\right) \frac{n-k}{k-1} = F_{1-\alpha, k-1, n-k}.$$

□

Treatment sum of squares: SSTR

Sample size: (Weights)	n_1	n_2	\dots	n_k	$n = \sum_{j=1}^k n_j$ <i>Weighted average</i>
Sample means:	$\bar{Y}_{.1}$	$\bar{Y}_{.2}$	\dots	$\bar{Y}_{.k}$	$\bar{Y}_{..} = \frac{1}{n} \sum_{j=1}^k n_j \bar{Y}_{.j}$
True means:	μ_1	μ_2	\dots	μ_k	$\mu = \frac{1}{n} \sum_{j=1}^k n_j \mu_j$
Squares:	$(\bar{Y}_{.1} - \bar{Y}_{..})^2$	$(\bar{Y}_{.2} - \bar{Y}_{..})^2$	\dots	$(\bar{Y}_{.k} - \bar{Y}_{..})^2$	SSTR

$$SSTR := \sum_{j=1}^k n_j (\bar{Y}_{.j} - \bar{Y}_{..})^2$$

1. When $k = 1$, $SSTR \equiv 0$.
2. When $k = 2$, say X_1, \dots, X_n and Y_1, \dots, Y_m :

$$\bar{Y}_{..} = \frac{1}{m+n} (n\bar{X} + m\bar{Y})$$

$$\begin{aligned} SSTR &= n \left[\bar{X} - \frac{1}{n+m} (n\bar{X} + m\bar{Y}) \right]^2 + m \left[\bar{Y} - \frac{1}{n+m} (n\bar{X} + m\bar{Y}) \right]^2 \\ &= n \left[\frac{m(\bar{X} - \bar{Y})}{n+m} \right]^2 + m \left[\frac{n(\bar{X} - \bar{Y})}{n+m} \right]^2 \\ &= \left[\frac{nm^2}{(n+m)^2} + \frac{n^2m}{(n+m)^2} \right] (\bar{X} - \bar{Y})^2 \\ &= \frac{nm}{n+m} (\bar{X} - \bar{Y})^2 \end{aligned}$$

$$SSTR = \frac{(\bar{X} - \bar{Y})^2}{\frac{1}{m} + \frac{1}{n}}$$

$$\begin{aligned}
SSTR &= \sum_{j=1}^k n_j (\bar{Y}_{.j} - \bar{Y}_{..})^2 = \sum_{j=1}^k n_j [(\bar{Y}_{.j} - \mu) - (\bar{Y}_{..} - \mu)]^2 \\
&= \sum_{j=1}^k n_j [(\bar{Y}_{.j} - \mu)^2 + (\bar{Y}_{..} - \mu)^2 - 2(\bar{Y}_{.j} - \mu)(\bar{Y}_{..} - \mu)] \\
&= \sum_{j=1}^k n_j (\bar{Y}_{.j} - \mu)^2 + \sum_{j=1}^k n_j (\bar{Y}_{..} - \mu)^2 - 2(\bar{Y}_{..} - \mu) \sum_{j=1}^k n_j (\bar{Y}_{.j} - \mu) \\
&= \sum_{j=1}^k n_j (\bar{Y}_{.j} - \mu)^2 + n(\bar{Y}_{..} - \mu)^2 - 2(\bar{Y}_{..} - \mu)n(\bar{Y}_{..} - \mu) \\
&= \sum_{j=1}^k n_j (\bar{Y}_{.j} - \mu)^2 - n(\bar{Y}_{..} - \mu)^2 \tag{12.2.1}
\end{aligned}$$

⇓

$$SSTR = \sum_{j=1}^k n_j \left[(\bar{Y}_{.j} - \mu_j)^2 - 2(\bar{Y}_{.j} - \mu_j)(\mu - \mu_j) + (\mu - \mu_j)^2 \right] - n(\bar{Y}_{..} - \mu)^2$$

Notice that

$$\bar{Y}_{.j} \sim N(\mu_j, \sigma^2/n_j) \quad \text{and} \quad \bar{Y}_{..} \sim N(\mu, \sigma^2/n)$$

\Rightarrow

$$\begin{aligned} \mathbb{E}[SSTR] &= \sum_{j=1}^k n_j \left[\frac{\sigma^2}{n_j} - 2 \times 0 + (\mu - \mu_j)^2 \right] - n \frac{\sigma^2}{n} \\ &= (k-1)\sigma^2 + \sum_{j=1}^k n_j (\mu - \mu_j)^2 \end{aligned}$$

Remark

When $\mu_1 = \dots = \mu_j$ then

0.1 $\mathbb{E}[SSTR] = (k-1)\sigma^2$

0.2 $MSTR := \frac{SSTR}{k-1}$ is an unbiased estimator for σ^2 .

0.3 $SSTR/\sigma^2 \sim \text{Chi square } (df = k-1)$.

(Homework)

Test $H_0 : \mu_1 = \cdots = \mu_k$ v.s. μ_j are not the same.

Case I. when σ^2 is known.

Reject H_0 if $SSTR/\sigma^2 \geq \chi_{1-\alpha, k-1}^2$.

Case II. when σ^2 is unknown.

.....

Sum of Squared Errors: SSE

1. Sum of squared error:

$$\begin{aligned}SSE &:= \sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{\cdot j})^2 \\&= \sum_{j=1}^k (n_j - 1) \left[\frac{1}{n_j - 1} \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{\cdot j})^2 \right] \\&= \sum_{j=1}^k (n_j - 1) S_j^2\end{aligned}$$

2. Pooled variance S_p^2 :

$$S_p^2 = \frac{SSE}{\sum_{j=1}^k (n_j - 1)} = \frac{SSE}{n - k}$$

Mean square of error $MSE = S_p^2$

Notice that

1. $(n_j - 1)S_j^2/\sigma^2 \sim \text{Chi square } (df = n_j - 1)$
2. S_j^2 's are independent
3. $SSE/\sigma^2 = (n - k)S_p^2/\sigma^2 = \sum_{j=1}^k (n_j - 1)S_j^2/\sigma^2$,
Sum of independent of Chi squares

↓

Thm. No matter $H_0 : \mu_1 = \dots = \mu_k$ is true or not

- a. $SSE/\sigma^2 = (n - k)S_p^2/\sigma^2 \sim \text{Chi square } (df = \sum_{j=1}^k (n_j - 1) = n - k)$
- b. $SSTR \perp SSE$.

Proof. We have shown part (a). Part (b) is trickier. Indeed, both parts are a consequence of **Cochran's theorem**¹ ... □

¹https://en.wikipedia.org/wiki/Cochran%27s_theorem

Let's see two special cases of

Thm. No matter $H_0 : \mu_1 = \dots = \mu_k$ is true or not

a. $SSE/\sigma^2 = (n - k)S_p^2/\sigma^2 \sim \text{Chi square } (df = \sum_{j=1}^k (n_j - 1) = n - k)$

b. $SSTR \perp SSE$.

Cases

1. $k = 1$, one sample case, S_p^2 is sample variance

Chapter 7

a. $(n - 1)S^2/\sigma^2 \sim \chi^2(n - 1)$

b. $SSTR \equiv 0$

2. $k = 2$, two sample case

Chapter 9

a. $(n - 2)S_p^2/\sigma^2 \sim \chi^2(n - 2)$

b. $\bar{X} - \bar{Y} \perp S_p^2 \iff SSTR \perp SSE$

Under $H_0 : \mu_1 = \dots = \mu_k$

1. $SSTR/\sigma^2 \sim \chi^2(k-1)$
2. $SSE/\sigma^2 \sim \chi^2(n-k)$
3. $SSTR \perp SSE$

$$\implies F = \frac{SSTR/(k-1)}{SSE/(n-k)} \sim F(df_1 = k-1, df_2 = n-k)$$

Reject H_0 if $F \geq F_{1-\alpha, k-1, n-k}$

Total Sum of Squares: SSTOT

$$SSTOT = SSE + SSTR$$

$$SSTOT := \sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{..})^2$$

||

$$\sum_{j=1}^k \sum_{i=1}^{n_j} \left[(Y_{ij} - \bar{Y}_{j.}) + (\bar{Y}_{j.} - \bar{Y}_{..}) \right]^2$$

||

$$\sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{j.})^2 + 2 \sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{j.}) (\bar{Y}_{j.} - \bar{Y}_{..}) + \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{Y}_{j.} - \bar{Y}_{..})^2$$

||

$$\sum_{j=1}^k \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{j.})^2 + 2 \sum_{j=1}^k (\bar{Y}_{j.} - \bar{Y}_{..}) \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{j.}) + \sum_{j=1}^k n_j (\bar{Y}_{j.} - \bar{Y}_{..})^2$$

||

$$SSE + 0 + SSTR$$

$$SSTOT = SSE + SSTR$$

↓

$$\frac{SSTOT}{\sigma^2} = \frac{SSE}{\sigma^2} + \frac{SSTR}{\sigma^2}$$

∝

∝

∝

$$\chi^2(n-1) \quad \chi^2(n-k) \quad \perp \quad \chi^2(k-1)$$

Under H_0

✓

Under H_0

One-way ANOVA Table

Source of Variance	Degree of Freedom (df)	Sum Square (SS)	Mean Square (MS)	F-ratio
Between Groups (Treatment)	k-1	$SSB = \sum_{j=1}^k \left(\frac{T_j^2}{n_j} \right) - \frac{T^2}{n}$ $SSB = \sum_{j=1}^k n_j (\bar{X}_j - \bar{X}_t)^2$	$MSB = \frac{SSB}{k-1}$	$F = \frac{MSB}{MSW}$
Within Groups (Error)	n-k	$SSW = \sum_{j=1}^k \sum_{i=1}^n X_{ij}^2 - \sum_{j=1}^k \left(\frac{T_j^2}{n_j} \right)$ $SSW = \sum_{j=1}^k \sum_{i=1}^n (X_{ij} - \bar{X}_j)^2$	$MSW = \frac{SSW}{n-k}$	
Total	n-1	$SST = \sum_{j=1}^k \sum_{i=1}^n X_{ij}^2 - \frac{T^2}{n}$ $SST = \sum_{j=1}^k \sum_{i=1}^n (X_{ij} - \bar{X}_t)^2$		

- $SST = SSB + SSW$

k: number of groups n: number of samples
df: degree of freedom

Source	df	SS	MS	F	P
Treatment	k - 1	SSTR	MSTR	$\frac{MSTR}{MSE}$	$P(F_{k-1, n-k} \geq \text{observed } F)$
Error	n - k	SSE	MSE		
Total	n - 1	SSTOT			

Common notation

d.f.

k-1 Error sum of squares
Mean square of error
(Pooled sample variance)

$$SSE = SSW = SS_{within}$$
$$MSE = MSW = MS_{within} = S_p^2$$

n-k Treatment sum of squares
Mean square of treatment

$$SSTR = SSB = SS_{between}$$
$$MSTR = MSB = MS_{between}$$

n-1 Total sum of squares:

$$SST = SSTOT$$

One way ANOVA v.s. Two sample t -test

Let X_1, \dots, X_n and Y_1, \dots, Y_m be samples from $N(\mu_X, \sigma^2)$ and $N(\mu_Y, \sigma^2)$, respectively.

Recall

$$1. SSTR/\sigma^2 = \frac{(\bar{X} - \bar{Y})^2}{\sigma^2 \left(\frac{1}{n} + \frac{1}{m}\right)} \sim \chi^2(1)$$

$$2. SSE/\sigma^2 = (n + m - 2)S_p^2/\sigma^2 \sim \chi^2(n + m - 2)$$

$$\Rightarrow F = \frac{SSTR/1}{SSE/(n + m - 2)} = \frac{(\bar{X} - \bar{Y})^2}{S_p^2 \left(\frac{1}{n} + \frac{1}{m}\right)} \sim F(df_1 = 1, df_2 = n + m - 2)$$

||
 T^2

$$\Rightarrow \alpha = \mathbb{P}(|T| \geq t_{\alpha/2, n+m-2}) = \mathbb{P}(T^2 \geq t_{\alpha/2, n+m-2}^2) = \mathbb{P}(F \geq F_{1-\alpha, 1, n+m-2})$$

Equivalent!

E.g. 1 Study the relation between smoking and heart rates.

Generations of athletes have been cautioned that cigarette smoking impedes performance. One measure of the truth of that warning is the effect of smoking on heart rate. In one study, six nonsmokers, six light smokers, six moderate smokers, and six heavy smokers each engaged in sustained physical exercise. Table 8.1.1 lists their heart rates after they had rested for three minutes.

	Nonsmokers	Light Smokers	Moderate Smokers	Heavy Smokers
	69	55	66	91
	52	60	81	72
	71	78	70	81
	58	58	77	67
	59	62	57	95
	65	66	79	84
<i>Averages:</i>	62.3	63.2	71.7	81.7

Show whether smoking affects heart rates at $\alpha = 0.05$.

Sol. Let μ_1, \dots, μ_4 be the true heart rates.

Test $H_0 : \mu_0 = \dots = \mu_4$ or not.

Critical region:

Let $\alpha = 0.05$. For these data, $k = 4$ and $n = 24$, so $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$ should be rejected if

$$F = \frac{SSTR/(4-1)}{SSE/(24-4)} \geq F_{1-0.05, 4-1, 24-4} = F_{.95, 3, 20} = 3.10$$

(see Figure 12.2.2).

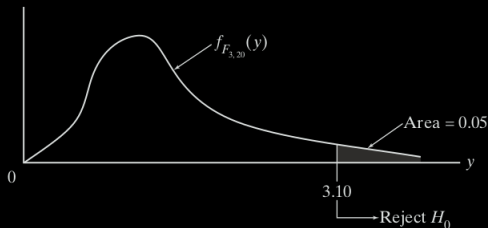


Figure 12.2.2

	Non smokers	Light Smokers	Moderate Smokers	Heavy Smokers
	69	55	66	91
	52	60	81	72
	71	78	70	81
	58	58	77	67
	59	62	57	95
	65	66	79	84
T_j	374	379	430	490
\bar{Y}_j	62.3	63.2	71.7	81.7

The overall sample mean, $\bar{Y}_{..}$, is given by

$$\begin{aligned}\bar{Y}_{..} &= \frac{1}{n} \sum_{j=1}^k T_j = \frac{374 + 379 + 430 + 490}{24} \\ &= 69.7\end{aligned}$$

Therefore,

$$\begin{aligned}SSR &= \sum_{j=1}^4 n_j (\bar{Y}_{.j} - \bar{Y}_{..})^2 = 6[(62.3 - 69.7)^2 + \dots + (81.7 - 69.7)^2] \\ &= 1464.125\end{aligned}$$

Similarly,

$$\begin{aligned}SSE &= \sum_{j=1}^4 \sum_{i=1}^6 (Y_{ij} - \bar{Y}_{.j})^2 = [(69 - 62.3)^2 + \dots + (65 - 62.3)^2] \\ &\quad + \dots + [(91 - 81.7)^2 + \dots + (84 - 81.7)^2] \\ &= 1594.833\end{aligned}$$

The observed test statistic, then, equals 6.12:

$$F = \frac{1464.125/(4 - 1)}{1594.833/(24 - 4)} = 6.12$$

Since $6.12 > F_{.95,3,20} = 3.10$, $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$ should be rejected. These data support the contention that smoking influences a person's heart rate.

Figure 12.2.3 shows the analysis of these data summarized in the ANOVA table format. Notice that the small P -value ($= 0.004$) is consistent with the conclusion that H_0 should be rejected.

Source	df	SS	MS	F	P
Treatment	3	1464.125	488.04	6.12	0.004
Error	20	1594.833	79.74		
Total	23	3058.958			

Figure 12.2.3



```
1 > Input <-c("
2 + rates group
3 + 69 non
4 + 52 non
5 + 71 non
6 + 58 non
7 + 59 non
8 + 65 non
9 + 55 light
10 + 60 light
11 + 78 light
12 + 58 light
13 + 62 light
14 + 66 light
15 + 66 moderate
16 + 81 moderate
17 + 70 moderate
18 + 77 moderate
19 + 57 moderate
20 + 79 moderate
21 + 91 heavy
22 + 72 heavy
23 + 81 heavy
24 + 67 heavy
25 + 95 heavy
26 + 84 heavy
27 + ")
28 > Data = read.table(textConnection(
      Input),
29 + header=TRUE)
```

```
1 > Data
2   rates  group
3 1    69   non
4 2    52   non
5 3    71   non
6 4    58   non
7 5    59   non
8 6    65   non
9 7    55  light
10 8    60  light
11 9    78  light
12 10   58  light
13 11   62  light
14 12   66  light
15 13   66 moderate
16 14   81 moderate
17 15   70 moderate
18 16   77 moderate
19 17   57 moderate
20 18   79 moderate
21 19   91  heavy
22 20   72  heavy
23 21   81  heavy
24 22   67  heavy
25 23   95  heavy
26 24   84  heavy
```

```

1 > # Check the levels
2 > levels(Data$group)
3 [1] "heavy" "light" "moderate" "non"
4 > # Order the groups
5 > Data$group <- ordered(Data$group, levels = c("non", "light", "moderate", "heavy")
6 )
7 > levels(Data$group)
8 [1] "non" "light" "moderate" "heavy"

```

```

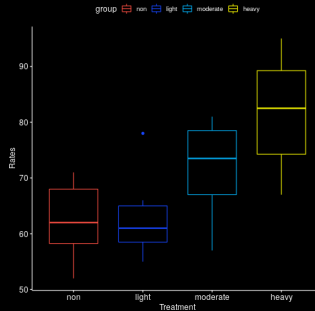
1 > # Compute summary statistics by groups
2 > # including count, mean, sd:
3 > library(dplyr) # a grammar of data manipulation
4 > group_by(Data, group) %>%
5 + summarise(
6 +   count = n(),
7 +   mean = mean(rates, na.rm = TRUE),
8 +   sd = sd(rates, na.rm = TRUE)
9 + )
10 # A tibble: 4 x 4
11   group   count mean  sd
12   <ord> <int> <dbl> <dbl>
13 1 non         6  62.3  7.26
14 2 light        6  63.2  8.16
15 3 moderate    6  71.7  9.16
16 4 heavy        6  81.7 10.8

```

```

1 # Box plots
2 # ++++++
3 # Plot rates by group and color by group
4 library(ggpubr)
5 png("Case_12-2-1-ggboxplot.png")
6 ggboxplot(Data, x = "group", y = "rates",
7           color = "group", palette = c("#00AFBB", "#E7B800", "#FC4E07", "blue"),
8           order = c("non", "light", "moderate", "heavy"),
9           ylab = "Rates", xlab = "Treatment")
10 dev.off()

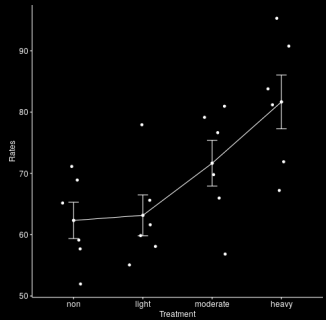
```




```

1 # Mean plots
2 # ++++++
3 # Plot rates by group
4 # Add error bars: mean_se
5 # (other values include: mean_sd, mean_ci, median_iqr, ....)
6 png("Case_12-2-1-ggline.png")
7 library(ggpubr)
8 ggline(Data, x = "group", y = "rates",
9         add = c("mean_se", "jitter"),
10        order = c("non", "light", "moderate", "heavy"),
11        ylab = "Rates", xlab = "Treatment")
12 dev.off()

```



```

1 > # Compute the analysis of variance
2 > res.aov <- aov(rates ~ group, data = Data)
3 > # Summary of the analysis
4 > summary(res.aov)
5           Df Sum Sq Mean Sq F value Pr(>F)
6 group         3  1464  488.0    6.12 0.00398 **
7 Residuals    20  1595   79.7
8 ---
9 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

1 > # Tukey multiple multiple-comparisons
2 > TukeyHSD(res.aov)
3   Tukey multiple comparisons of means
4     95% family-wise confidence level
5
6 Fit: aov(formula = rates ~ group, data = Data)
7
8 $group
9           diff          lwr          upr          p adj
10 light-non    0.8333333 -13.596955  15.26362  0.9984448
11 moderate-non 9.3333333  -5.096955  23.76362  0.2978123
12 heavy-non   19.3333333  4.903045  33.76362  0.0063659
13 moderate-light 8.5000000  -5.930289  22.93029  0.3755571
14 heavy-light  18.5000000  4.069711  32.93029  0.0091463
15 heavy-moderate 10.0000000  -4.430289  24.43029  0.2438158

```

1. diff: difference between means of the two groups
2. lwr, upr: the lower and the upper end point of the C.I. at 95% (default)
3. p adj: p-value after adjustment for the multiple comparisons

Inferences

if p-value \leq 0.05 \iff if zero is in the C.I.

```

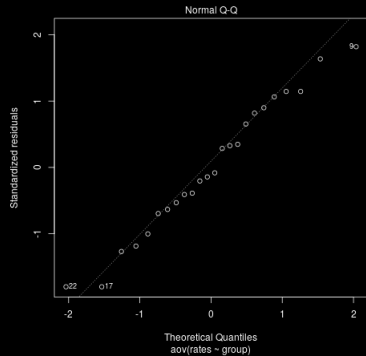
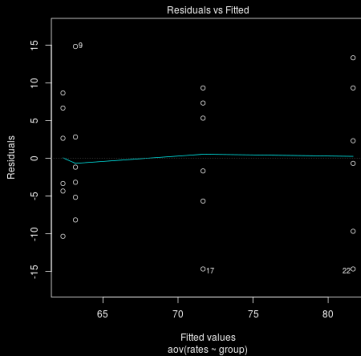
1 > # Or one may use multcomp package or multiple comparisons
2 > library(multcomp)
3 > summary(glht(res.aov, linfct = mcp(group = "Tukey")))
4
5 Simultaneous Tests for General Linear Hypotheses
6
7 Multiple Comparisons of Means: Tukey Contrasts
8
9
10 Fit: aov(formula = rates ~ group, data = Data)
11
12 Linear Hypotheses:
13             Estimate Std. Error t value Pr(>|t|)
14 light - non == 0    0.8333    5.1556  0.162 0.99844
15 moderate - non == 0  9.3333    5.1556  1.810 0.29776
16 heavy - non == 0   19.3333    5.1556  3.750 0.00629 **
17 moderate - light == 0 8.5000    5.1556  1.649 0.37544
18 heavy - light == 0  18.5000    5.1556  3.588 0.00901 **
19 heavy - moderate == 0 10.0000    5.1556  1.940 0.24382
20 ---
21 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
22 (Adjusted p values reported -- single-step method)

```

```

1 # Check ANOVA assumptions: test validity?
2 # diagnostic plots
3 layout(matrix(c(1,2),1,2)) # optional 1x2 graphs/page
4 plot(res.aov,c(1,2))

```



1. Residuals vs Fitted: test homogeneity of variances

One can also use Levene's test for this purpose:

```
1 > # Use Levene's test to test homogeneity of variances
2 > library(car)
3 > leveneTest(rates ~ group, data = Data)
4 Levene's Test for Homogeneity of Variance (center = median
5   )
6   Df F value Pr(>F)
7   group 3  0.3885 0.7625
8   20
```

2. Normal Q-Q plot: Test normality. (It should be close to diagonal line.)

One can also use Shapiro-Wilk test:

```
1 # Extract the residuals
2 > aov_residuals <- residuals(object = res.aov )
3 > # Run Shapiro-Wilk test
4 > shapiro.test(x = aov_residuals )
5
6 Shapiro-Wilk normality test
7
8 data:  aov_residuals
9 W = 0.9741, p-value = 0.7677
```

Non-parametric alternative to one-way ANOVA test

```
1 > # Non-parametric alternative to one-way ANOVA test
2 > # a non-parametric alternative to one-way ANOVA
3 > # is Kruskal-Wallis rank sum test, which can be
4 > # used when ANNOVA assumptions are not met.
5 > kruskal.test(rates ~ group, data = Data)
6
7   Kruskal-Wallis rank sum test
8
9   data: rates by group
10  Kruskal-Wallis chi-squared = 10.729, df = 3, p-value =
    0.01329
```

See Section 4 of Chapter 14 for more details.