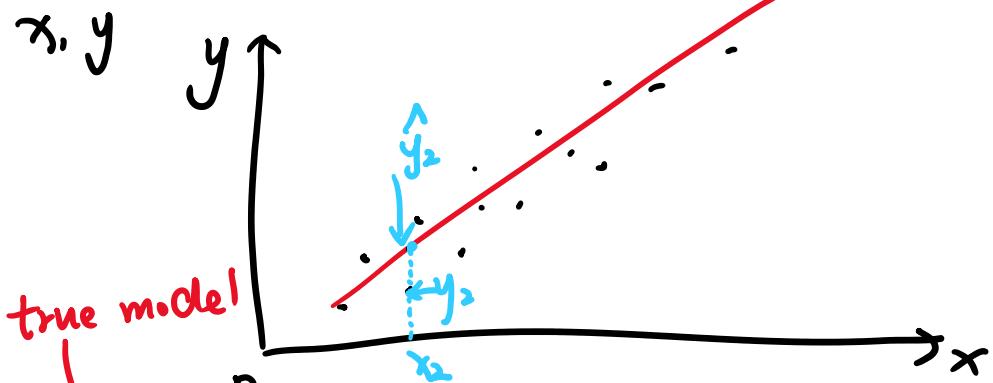


# Chapter 11: Simple Linear Regression.



$y = \beta_0 + \beta_1 x$ ,  $\Gamma$  correlation coefficient.

Suppose  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$\hat{y} = b_0 + b_1 x, \quad \hat{y} \Rightarrow y \\ b_0 \Rightarrow \beta_0 \\ b_1 \Rightarrow \beta_1$$

Residuals:

$$\hat{e}_i = y_i - \hat{y}_i$$

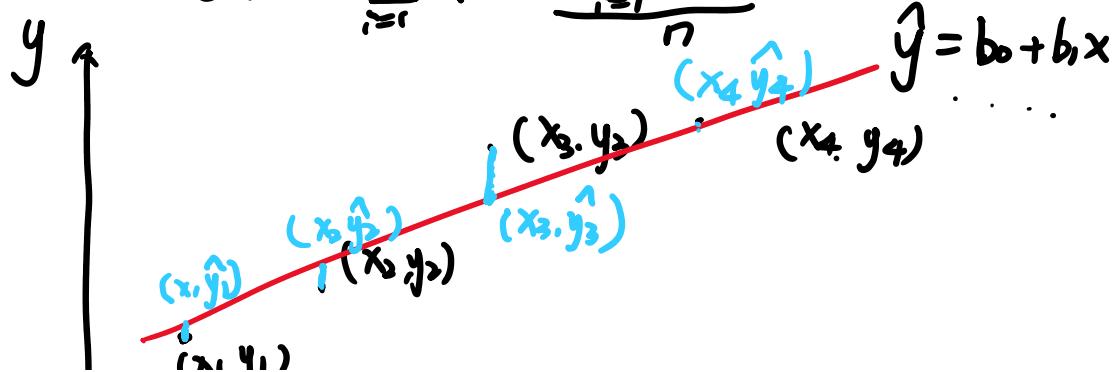
Least square method.  $(\bar{x} = \hat{\mu})$

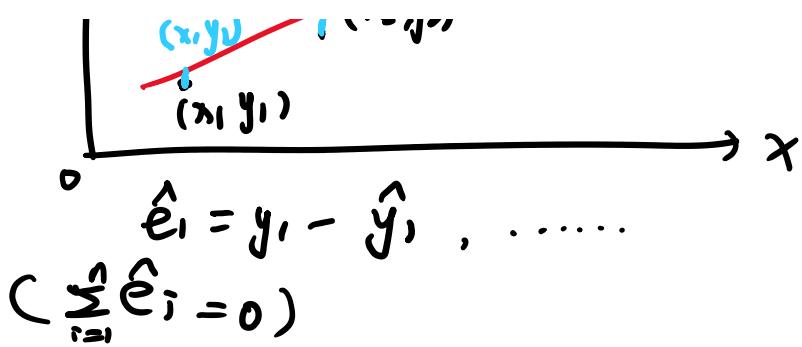
$$L = \sum_{i=1}^n \hat{e}_i^2 = \sum (y_i - \hat{y}_i)^2 \leftarrow \text{minimize } L.$$

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned} \quad \Rightarrow \quad \begin{cases} b_1 = \frac{S_{xy}}{S_{xx}} = \hat{\beta}_1 \\ b_0 = \bar{y} - b_1 \bar{x} = \hat{\beta}_0 \end{cases}$$

$$S_{xy} = \sum_{i=1}^n x_i y_i - \frac{(\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n}$$

$$S_{xx} = \sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}$$





Example:

$x$	1	2	4	7
$y$	5	3	2	1

Find regression line.

$$S_{xx} = \sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}$$

$$= 70 - \frac{14^2}{4} = 21$$

$$S_{xy} = \sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n} = 26 - \frac{14 \cdot 11}{4} = -\frac{25}{2}$$

$$b_1 = \hat{\beta}_1 = -\frac{25}{42}$$

$$b_0 = \hat{\beta}_0 = \frac{29}{6}$$

$$\Rightarrow \hat{y} = \frac{29}{6} - \frac{25}{42} x \quad z = x^2$$

$$y = ax^2 + b$$

$$\Rightarrow y = az + b$$

transformation  
for non-linear.

What's the residual for  $(x_i, y_i)$

$$x_1 = 1, y_1 = 5$$

$$\hat{y}_1 = \frac{29}{6} - \frac{25}{42} = \frac{29 \cdot 7 - 25}{42} = 4.238$$

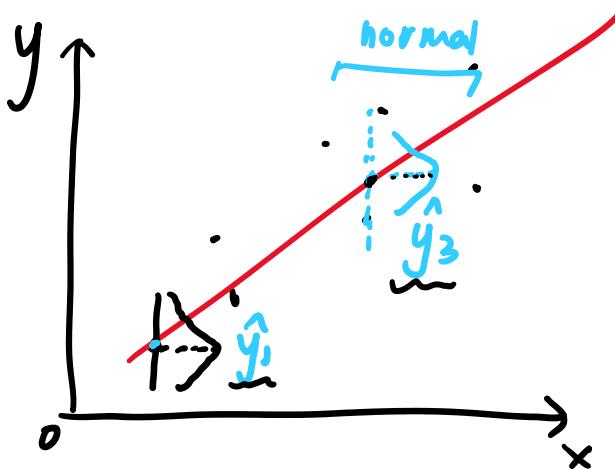
$$y_1 = 5$$

$$\hat{e}_1 = y_1 - \hat{y}_1 = 5 - 4.238 = 0.762$$

$$\epsilon_i = y_i - \hat{y}_i \rightarrow \text{real error}$$

Assumption:

1.  $E(\epsilon_i) = 0$
2.  $V(\epsilon_i) = \sigma^2 \leftarrow \text{SSG}$
3.  $\text{Cov}(\epsilon_i, \epsilon_j) = 0$
4. Normal for  $\epsilon_i$



Estimate  $\sigma^2$ .

The sum of square error

$$SSE = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\hat{\sigma}^2 = \frac{SSE}{n-2} \Rightarrow [E(\hat{\sigma}^2) = \sigma^2] \text{ unbiased}$$

The total sum of squares.

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - n\bar{y}^2$$

The regression sum of squares:

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SST = SSE + SSR$$

Properties of Least squares estimator:

$$(1) E(\hat{\beta}_1) = \beta_1 \quad \text{or} \quad E(b_1) = \beta_1$$

$$(2) V(\hat{\beta}_1) = V(b_1) = \frac{\sigma^2}{S_{xx}} \quad \begin{matrix} \text{variance for residuals} \\ \text{or errors} \end{matrix}$$

(3)  $\hat{\beta}_1$  is normal R.V

as in our assumption,  $\epsilon \sim N(0, \sigma^2)$

As in our assumption,  $\epsilon \sim N(0, \sigma^2)$

$$\hat{y} = b_0 + b_1 x + \epsilon$$

$\hat{y}$  Normal  
 $\epsilon$  Normal

Note: If  $b_1 = 0$  then,  $\hat{y}$  has no linear relation with  $x$

### Hypothesis test for simple linear regression.

Test for  $\beta_1 = 0$ .

1. Hypothesis:  $H_0: \beta_1 = 0$

$H_A: \beta_1 \neq 0$  (no linear relationship)

2.  $T = \frac{\hat{\beta}_1}{\hat{\sigma}/\sqrt{s_{xx}}}$  where  $\hat{\sigma} = \sqrt{\frac{SSE}{n-2}}$

Degree of freedom =  $n - 2$ .

3. Critical value.

4. Conclusion.

Example:

$H_0: \beta_1 = 0$

( $n=4$ )

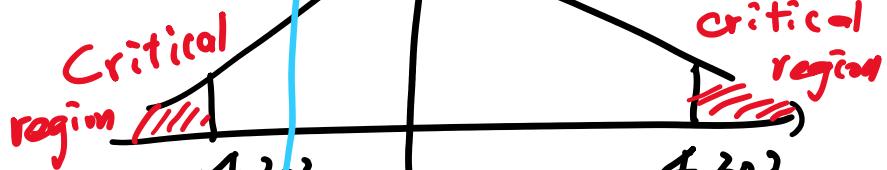
$H_A: \beta_1 \neq 0$

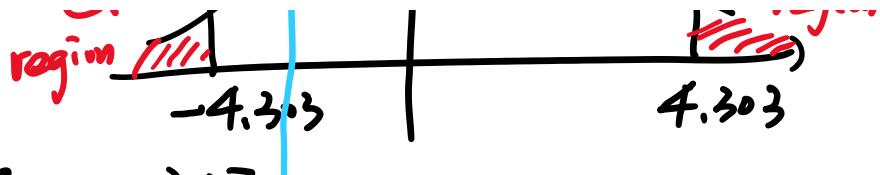
$$T = \frac{\hat{\beta}_1}{\hat{\sigma}/\sqrt{s_{xx}}} = \frac{-25/42}{\sqrt{\frac{55/42}{4-2}}/\sqrt{2}} = -3.371$$

$\alpha = 5\%$

$t_{2, \frac{5}{2}\%} = 4.303$

-3.371  $\leq$  test statistic





$\therefore \text{As } t_2, \frac{5}{2}\% < -3.371$

do not reject  $H_0$

## ANOVA: Analysis of Variance.

Recall:

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SST = SSR + SSE.$$

$$F = \frac{SSR}{SSE/(n-2)} \sim F \text{ distribution with d.f. } (1, n-2)$$

### Hypothesis Test (ANOVA)

1. Hypothesis:  $H_0: \beta_1 = 0$

$$H_A: \beta_1 \neq 0$$

### 2. Test statistic

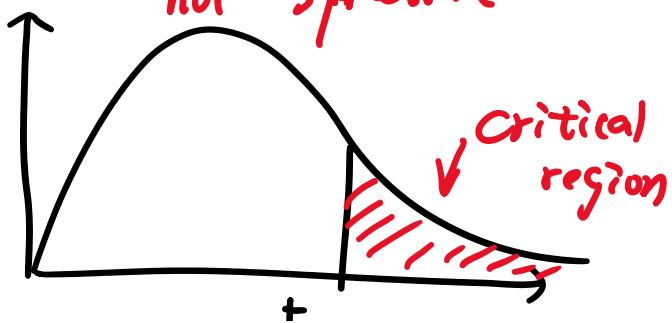
$$F = \frac{SSR}{SSE/(n-2)} \sim F(1, n-2)$$



always one tail as  $F$  is a positive  
not symmetric distribution.

### 3. Critical value

$$F_{1, n-2, \alpha}$$



$$F_{1, n-2, \alpha}$$

$$F \xrightarrow{\text{---}} F_{1, n-2, \alpha} \quad \text{---} \quad \text{---}$$

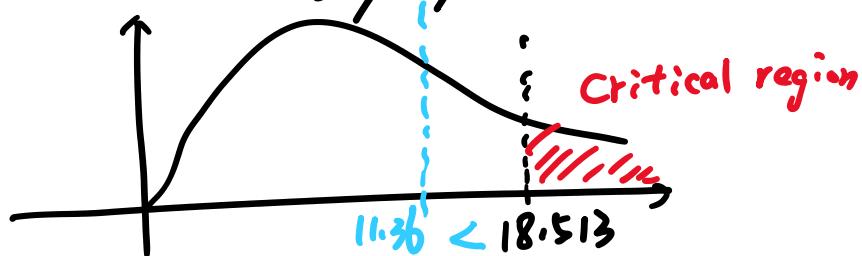
#### 4. Conclusion

If  $F > F_{1, n-2, \alpha}$ , then reject  $H_0$   
 otherwise - do not reject  $H_0$ .

Example:  $n=4$ ,  $\alpha = 5\%$

critical value  $\rightarrow F_{1, 2, 5\%} = 18.513$

$$\therefore F = \frac{(-25/42)(-25/2)}{55/42/2} = 11.36$$



$\therefore$  do not reject  $H_0$

ANOVA Table:

Source of Variation	Sum of squares	Degree of Freedom	Mean Square	F
Regression	SSR	1	$MSR = \frac{SSR}{1}$	$\frac{MSR}{MSE}$
Error	SSE	$n-2$	$MSE = \frac{SSE}{n-2}$	
Total	SST	$n-1$		

Example

	SS	D.F	MS	F
Regression	$625/84$	1	$625/84$	$625/55 = 11.36$
Error	$55/42$	$4-2=2$	$55/84$	
	$25/1$			

Error	<del>55/42</del>	<del>7/6</del>	<del>-7/6</del>
Total	<del>25/4</del>		

Confidence Intervals:

1. Slope  $\beta_1$ :

$$\hat{\beta}_1 + t_{n-2, \frac{\alpha}{2}} \cdot \frac{\hat{\sigma}}{\sqrt{S_{xx}}} \quad \text{circled } \hat{\sigma} \quad \leftarrow SD(\hat{\beta}_1)$$

2. mean of  $y$  ( $\bar{y}$ ) at  $x=x_0$

$$\hat{y}_0 \pm t_{n-2, \frac{\alpha}{2}} \cdot \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}$$

where  $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$

3. Predicting new values. (prediction interval)

for a single  $y_0$  at  $x=x_0$

$$\hat{y}_0 \pm t_{n-2, \frac{\alpha}{2}} \cdot \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}$$

where  $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$

Example:

1. A 95% C.I. for  $\beta_1$ :  $\hat{\sigma}$

$$-\frac{25}{42} \pm t_{2, 0.025} \cdot \frac{\sqrt{55/84}}{\sqrt{21}} \leftarrow S_{xx}$$

$$= (-1.355, 0.646)$$

2. A 95% C.I. for  $\bar{y}$  at  $x_0=5$

2. A 95% CI for  $\hat{Y}$  at  $x_0=5$

$$\left( \frac{29}{5} - \frac{25}{42} \cdot 5 \right) \pm t_{2,0.025} \cdot \sqrt{\frac{s^2}{8}} \left[ \frac{1}{4} + \frac{(5 - \frac{14}{4})^2}{s^2} \right]$$

3. A 95% CI for  $y_0$  at  $x_0=5$

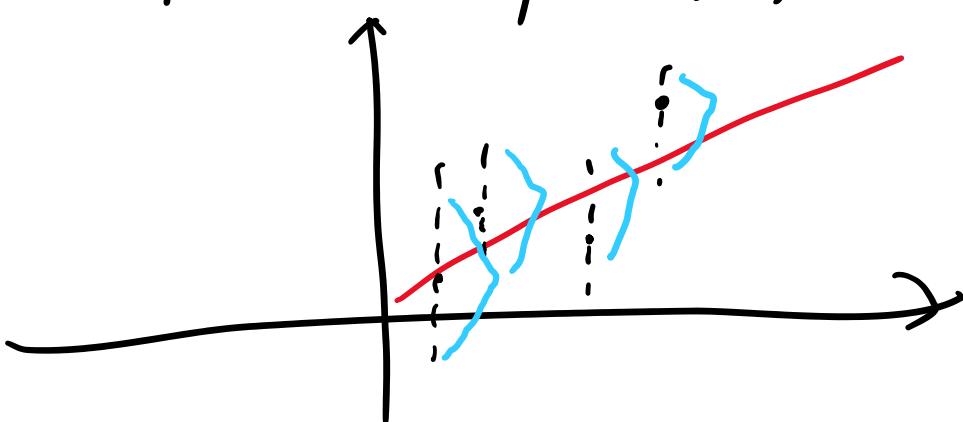
$$(-2.196, 6.307)$$

$$\sqrt{1 - \frac{1}{\text{circled value}}}$$

### Assumptions of Regression Model:

Example:

$x$	$y$	$\hat{y}_i$	$e_i$
1	5	4.238	0.762
2	3	3.6429	-0.6429
4	2	2.4524	-0.4524
7	1	0.6667	0.3333



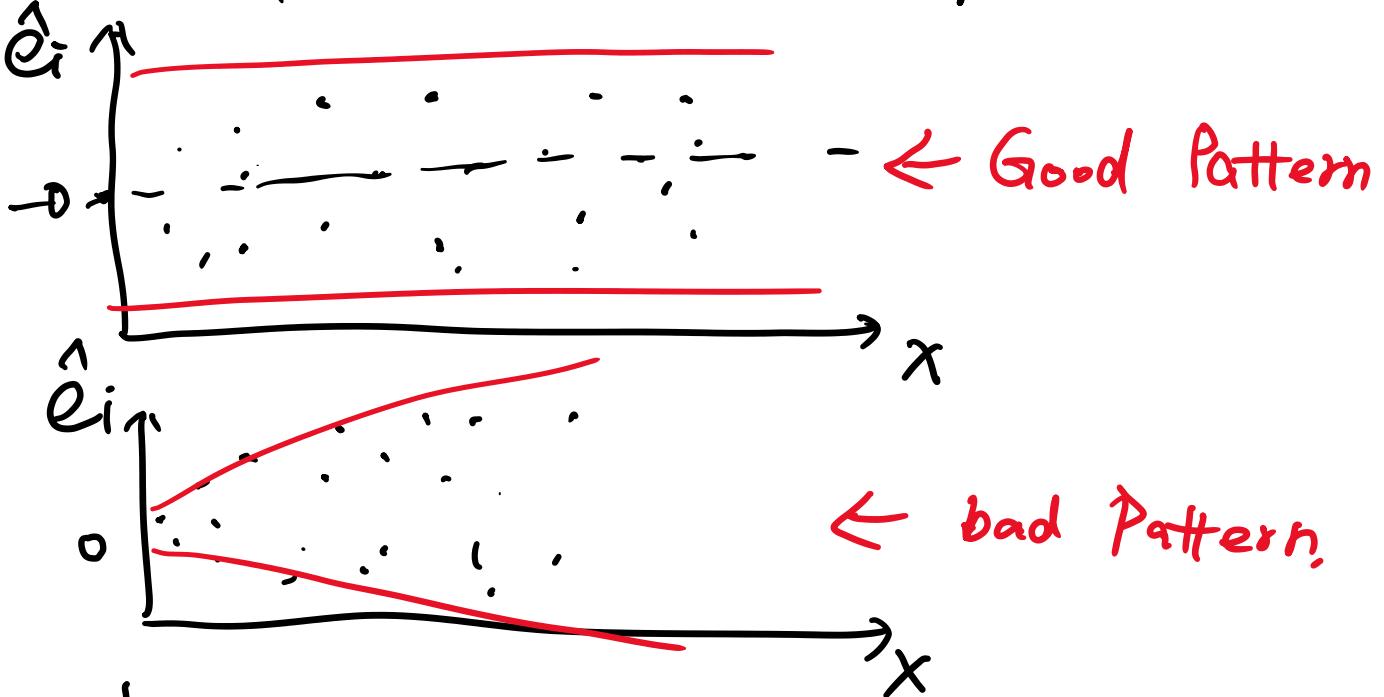
$$V(\varepsilon) = \sigma^2 \Leftarrow \text{Equal Variances}$$

1. Check Normal Assumptions

to check if  $\varepsilon_i \sim N(0, \sigma^2)$ , do a normal probability plot of the observed residuals

2. Check Equal Variance Assumptions.

## 2 Check Equal Variance Assumptions.



Recall.  $R$  ← correlation coefficient.

$R^2$  ← coefficient of determination.

$$\begin{aligned}
 r = R &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(n-1)S_x S_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \\
 &= \frac{S_{xy}}{S_x \times SST} \\
 &= \sqrt{\frac{SSR}{SST}}
 \end{aligned}$$

$$\therefore R^2 = \frac{SSR}{SST}$$

$R$  is an estimation of the population coefficient of  $\rho$ .

$$R = r = \hat{\rho}$$

$$R = r = p$$

## Chapter 13: Analysis of Variance.

### Overall F-Test

Example:

$n_1=7$	$n_2=5$	$n_3=7$	$n_4=8$
7.5	5.8	.5.9	6.2
6.2	7.3	6.2	6.8
6.9	8.2	5.8	5.7
7.4	7.1	4.7	4.9
9.2	7.8	8.3	6.2
8.5		7.2	7.1
7.6		6.2	5.4

ANOVA to test mean of groups are equal.

$$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$$

$H_A$ : At least one of them not equal.

Notation:

$$a = \# \text{ of groups / treatment} \quad a=4$$

$$\bar{y}_{i\cdot} = \text{total of the } i\text{th group.} \quad \bar{y}_{1\cdot} = 7.5 + \dots + 7.6$$

$$\bar{\bar{y}}_{i\cdot} = \text{mean of the } i\text{th group.} \quad \bar{\bar{y}}_{1\cdot} = \frac{\bar{y}_{1\cdot}}{7} = 7.59$$

$$\underline{\bar{y}\ldots} = \text{grand total}$$

$$\frac{y_{..} = \bar{y}}{\bar{y}_{..} = \text{grand mean}}$$

$N$  = total number observations.

$n_i$  = # of observations in the  $i$ th group.

$y_{ij}$  =  $i$ th group's  $j$ th observations.

$s_i$  = sample standard deviation for the  $i$ th group.

Example:

$$y_{1.} = 53.1 \quad y_{2.} = 36.2 \quad y_{3.} = 44.3 \quad y_{4.} = 48.1$$

$$\bar{y}_{1.} = 7.59 \quad \bar{y}_{2.} = 7.24 \quad \bar{y}_{3.} = 6.33 \quad \bar{y}_{4.} = 6.01$$

$$n_1 = 7 \quad n_2 = 5 \quad n_3 = 7 \quad n_4 = 8$$

$$y_{..} = 181.7$$

$$\bar{y}_{..} = 6.73$$

$$N = 27$$

Hypothesis:  $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$   
 $H_A: \text{At least one of them not equal}$

Test statistic:

$$F_{a-1, N-a} = \frac{SST_F / (a-1)}{SSE / (N-a)}$$

where:  $SST = \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^a \sum_{j=1}^{n_i} y_{ij}^2 - \frac{y_{..}^2}{N}$

variance of

$i=1 \quad j=1 \quad \dots \quad i=1 \quad j=N \quad N$

Variance of total

$$d.f = N-1$$

$$SST_r = \sum_{i=1}^q n_i (\bar{y}_{i\cdot} - \bar{y}_{..})^2 = \sum_{i=1}^q \frac{\bar{y}_{i\cdot}^2}{n_i} - \frac{\bar{y}_{..}^2}{N}$$

Similar to SSR

Variance among groups

$$d.f = q-1$$

$$SSE = \sum_{i=1}^q \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{ij})^2 = \sum_{i=1}^q (n_i - 1) S_i^2$$

Variance within groups

$$d.f = N-q$$

ANOVA TABLE:

Source of Variance	SS	D.F.	MS.	F
Treatment	$SST_r$	$q-1$	$\frac{SST_r}{q-1} = MST_r$	$\frac{MST_r}{MSE}$
Error	$SSE$	$N-q$	$\frac{SSE}{N-q} = MSE$	
Total	$SST$	$N-1$		

Example :

	SS	D.F.	MS	F
Treatments	11.673 $SST_r$	$3 \quad q-1$	3.891	$\frac{3.891}{0.8828} = 4.4$
Error	20.3043 $SSE$	$23 \quad N-q$	0.8828	
Total	31.9773 $SST$	$26 \quad N-1$		

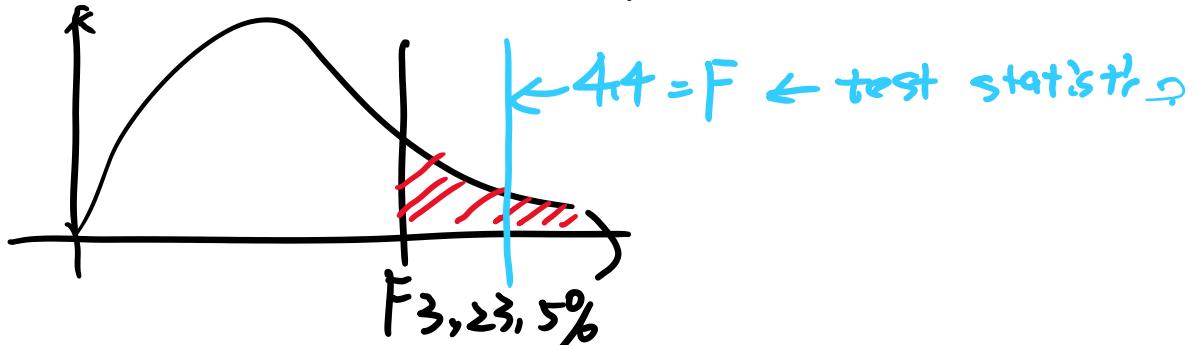
Critical value

$$F_{(q-1, N-q)} \alpha \% = F_{3, 23, 5\%}$$

$$F_{(a-1), (N-a), 5\%} = F_{3, 23, 5\%}$$

$$= 3.03$$

$$F = 4.4 > F_{3, 23, 5\%} = 3.03$$



$\therefore$  we reject  $H_0$ . conclude that at least one of the mean is not equal with others.

Test for individual pairs of means.

Fisher's LSD (Least Significant Difference)

1. Hypothesis:  $H_0: \mu_i = \mu_j$

$H_A: \mu_i \neq \mu_j$ ,  $\forall i, j$  any pair.

$$\binom{4}{2} = 6$$

2. Test statistic.

$$\textcircled{1} \quad \bar{y}_{ij} - \bar{y}_{jk}$$

$$\textcircled{2} \quad \text{LSD} = t_{n-a, \frac{\alpha}{2}} \cdot \sqrt{MSE \left( \frac{1}{n_i} + \frac{1}{n_j} \right)}$$

3. Critical Region:

Reject  $H_0$  if

$$| \bar{y}_{ij} - \bar{y}_{jk} | > \text{LSD}$$

$|\bar{y}_{i\cdot} - \bar{y}_{j\cdot}| > LSD$   
 otherwise do not reject  $H_0$ .

#### 4. Conclusion.

Example: Test  $\mu_1$  and  $\mu_2$  whether have a difference at  $\alpha = 5\%$ .

$$\text{Hypothesis: } H_0: \mu_1 = \mu_2$$

$$H_A: \mu_1 \neq \mu_2.$$

Test statistic:

$$\begin{aligned} |\bar{y}_{1\cdot} - \bar{y}_{2\cdot}| &= |7.59 - 7.24| = 0.35 \\ LSD &= t_{23, 0.025} \sqrt{MSE \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} \\ &= 2.059 \cdot \sqrt{0.8828 \left( \frac{1}{7} + \frac{1}{5} \right)} \\ &= 1.138 \end{aligned}$$

$\therefore 0.35 < LSD$   
 $\therefore$  do not reject  $H_0$ .

Hypothesis	LSD	$ \bar{y}_{i\cdot} - \bar{y}_{j\cdot} $	Conclusion
$H_0: \mu_1 = \mu_2$	1.138	0.35	
$\mu_1 = \mu_3$	1.039	1.26	$\mu_1 \neq \mu_3$
$\mu_1 = \mu_4$	1.006	1.58	$\mu_1 \neq \mu_4$

$H_0: \mu_1 = \mu_2$	1.006	0.58	$H_0: \mu_1 = \mu_2$
$H_0: \mu_2 = \mu_3$	1.138	0.91	
$H_0: \mu_2 = \mu_4$	1.108	1.23	$H_0: \mu_2 \neq \mu_4$
$H_0: \mu_3 = \mu_4$	1.061	0.32	

Assumptions for ANOVA:

1. Each population must be normal.
2. The populations have equal variances.
3. The samples must be independent.
4. All other factors must be equal other than groups.

Residuals for ANOVA:

$$e_{ij}^1 = y_{ij} - \bar{y}_{\cdot i} \quad (\text{observation minus sample average})$$

Assumption: 1. residuals follow a normal distribution.

2. For equal group, the variance of residuals should be similar.