10.1 : Inference about difference between two population means , $\sigma_1, \sigma_2$ Known .

→ Assumptions

1. Sample 1 Random

2. Sample 2 Random

3. Sample 1 and sample 2 independent .

4. a. population 1 Normal .

   b. population 2 Normal .

OR Both sample are large enough .
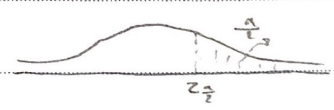
★ large enough samples : $n_1 \geq 30$  $\Big\}$ use it if pop 1,2 not Normal .
   $n_2 \geq 30$

→ point estimator for $\mu_1 - \mu_2 = \bar{X}_1 - \bar{X}_2$ .

→ confidence$^{CI}$ interval / interval estimate for $\mu_1 - \mu_2 = (\bar{X}_1 - \bar{X}_2) \mp E$

   margin of error $(E) = Z_{\frac{\alpha}{2}} \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}$   standard error .

→ $(1-\alpha)$ CI for $\mu_1 - \mu_2 = (\bar{X}_1 - \bar{X}_2) \mp Z_{\frac{\alpha}{2}} \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}$ .



→ Hypothesis test about $\mu_1 - \mu_2$ :

① lower tail test

$H_0 : \mu_1 - \mu_2 \geq D_0$

$H_1 : \mu_1 - \mu_2 < D_0$

② upper tail test

$H_0 : \mu_1 - \mu_2 \leq D_0$

$H_1 : \mu_1 - \mu_2 > D_0$

③ Two tail test .

$H_0 : \mu_1 - \mu_2 = D_0$

$H_1 : \mu_1 - \mu_2 \neq D_0$ .
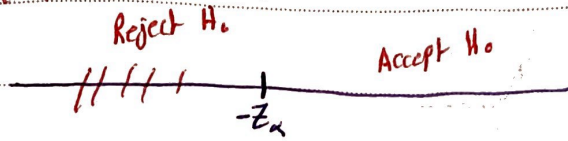
→ Test statistic :

$$Z = \dfrac{(\bar{X}_1 - \bar{X}_2) - D_0}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}}$$
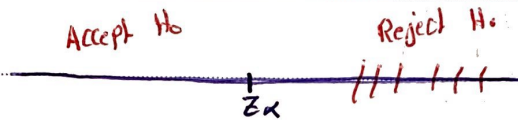
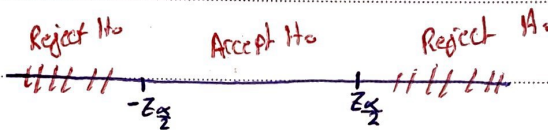By Z-table .

$\rightarrow$ Critical value Approch :

① Lower tail test :


Reject $H_0$     Accept $H_0$     $-Z_\alpha$

Reject $H_0$ if $Z \leq -Z_\alpha$ .

② Upper tail test :


Accept $H_0$     Reject $H_0$     $Z_\alpha$

Reject $H_0$ if $Z \geq Z_\alpha$ .

③ Two tail test :


Reject $H_0$     Accept $H_0$     Reject $H_0$     $-Z_{\frac{\alpha}{2}}$     $Z_{\frac{\alpha}{2}}$

Reject $H_0$ if $Z \geq Z_{\frac{\alpha}{2}}$ or $Z \leq -Z_{\frac{\alpha}{2}}$ .

$\rightarrow$ P-Value Approch :

Reject $H_0$ if P-Value $\leq \alpha$

① lower tail test :


p-value     Z

② upper tail test :


p-value     Z

③ Two tail test :


p-value     p-value

**6.2:** Inferences about $\mu_1 - \mu_2$ , $\delta_1$ and $\delta_2$ unknown.

→ **Assumptions:**

1. sample 1 Random sample from pop. 1

2. sample 2 Random sample from pop. 2

3. sample 1 and sample 2 are independent.

4. pop. 1 and pop. 2 have Normal distribution OR sample 1 and sample 2 are large enough.

$*$ large enough : $n_1 + n_2 \geq 20$ st $n_1 \approx n_2$

→ point estimater for $\mu_1 - \mu_2 = \bar{x}_1 - \bar{x}_2$

→ standard error $= \sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}$

→ margen of error $(E) = t_{\frac{\alpha}{2}} \sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}$

→ $df = \left\lfloor \dfrac{\left(\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}\right)^2}{\left(\dfrac{1}{n_1-1}\right)\left(\dfrac{s_1^2}{n_1}\right)^2 + \left(\dfrac{1}{n_2-1}\right)\left(\dfrac{s_2^2}{n_2}\right)^2} \right\rfloor$  →  $\lfloor 9.4 \rfloor = 9$

→ $1-\alpha$ CI $= (\bar{x}_1 - \bar{x}_2) \mp E$ .

→ **Hypotheses testing for $\mu_1 - \mu_2$ :**

① lower tail test

$H_0 : \mu_1 - \mu_2 \geq D_0$

$H_1 : \mu_1 - \mu_2 < D_0$

② upper tail test

$H_0 : \mu_1 - \mu_2 \leq D_0$

$H_1 : \mu_1 - \mu_2 > D_0$

③ two tail test

$H_0 : \mu_1 - \mu_2 = D_0$

$H_1 : \mu_1 - \mu_2 \neq D_0$

→ **Test statistic:**

$t = \dfrac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}}$

→ Critical Value Approach:

① lower tail test:

We reject H₀ if $t_{test} \leq -t_\alpha$.

② upper tail test:

We reject H₀ if $t_{test} \geq t_\alpha$.

③ two tail test:

We reject H₀ if $t_{test} \geq t_{\frac{\alpha}{2}}$ or $t_{test} \leq -t_{\frac{\alpha}{2}}$.

→ P-value Approach:

Reject H₀ if $P\text{-value} \leq \alpha$.

① lower tail test:



② upper tail test:



③ two tail test:



→ If we additionally have $\delta_1 = \delta_2$:

✱ Test statistic:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{s\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

✱ $df = n_1 + n_2 - 2$

✱ Pooled sample variance:

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

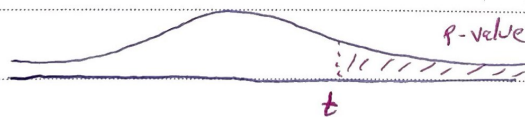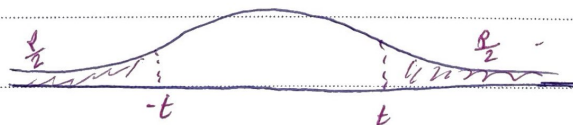b.3 : Inference about the difference between two population means, matched samples.

→ $H_0 : \mu_d = \mu_{d_0}$

$H_1 : \mu_d \neq \mu_{d_0}$ } two tailed test

- $n$ : sample size / # of element
- $\mu_1$ : pop.1 mean
- $\mu_2$ : pop 2 mean
- $\mu_d = \mu_1 - \mu_2$

→ Test statistic :

$$t = \frac{\bar{d} - \mu_{d_0}}{\frac{S_d}{\sqrt{n}}} \quad , \quad df = n-1$$

- $d_i = X_i^{I} - X_i^{II}$

- $\bar{d} = \frac{\sum d_i}{n}$

- $S_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}}$

→ Reject $H_0$ if $p$-value $\leq \alpha$

$p$-value = area in both tails.

Reject $H_0$ if $|t| \geq t_{\frac{\alpha}{2}}$ , $df = n-1$.

→ $(1-\alpha)$ CI for $\mu_d = \bar{d} \mp \boxed{t_{\frac{\alpha}{2}} \frac{S_d}{\sqrt{n}}}$ → margin of error

standard error .

10.4: Inferences about the difference between two population proportions.

→ Assumptions:

1. sample 1 and sample 2 Random.

2. samples 1 and 2 are independent.

3. samples 1 and 2 large enough.

    * large enough:

    pop : $n_1 \pi_1 \geq 5$ , $n_1(1-\pi_1) \geq 5$

              $n_2 \pi_2 \geq 5$ , $n_2(1-\pi_2) \geq 5$

* Notations:

$\pi_1$: proportion in pop. 1

$\pi_2$: proportion in pop. 2

$p_1$: proportion in sample 1

$p_2$: proportion in sample 2

$n_1$: sample 1 size.

$n_2$: sample 2 size.

    sample : $n_1 p_1 \geq 5$ , $n_1(1-p_1) \geq 5$

            $n_2 p_2 \geq 5$ , $n_2(1-p_2) \geq 5$

→ point estimator for $\pi_1 - \pi_2 = p_1 - p_2$ .

→ $(1-\alpha)$ CI for $\pi_1 - \pi_2 = (p_1 - p_2) \mp E$ .

→ margin of error $(E) = Z_{\frac{\alpha}{2}} \sqrt{\dfrac{p_1(1-p_1)}{n_1} + \dfrac{p_2(1-p_2)}{n_2}}$

→ standard error of $p_1 - p_2$ $(\delta_{p_1-p_2}) = \sqrt{\dfrac{\pi_1(1-\pi_1)}{n_1} + \dfrac{\pi_2(1-\pi_2)}{n_2}}$.

→ Hypotheses test about $\pi_1 - \pi_2$ :

① lower tail test :

$H_0 : \pi_1 - \pi_2 \geq 0$

$H_1 : \pi_1 - \pi_2 < 0$

② upper tail test:

$H_0 : \pi_1 - \pi_2 \leq 0$

$H_1 : \pi_1 - \pi_2 > 0$

③ two tail test :

$H_0 : \pi_1 - \pi_2 = 0$

$H_1 : \pi_1 - \pi_2 \neq 0$ .

→ **Remark:** under the $H_0$ when $H_0$ is true an equality we get $\pi_1 = \pi_2 = \pi$.

    → standard error $(\pi_1 = \pi_2 = \pi)$ : $\delta_{P_1 - P_2} = \sqrt{\pi(1-\pi)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$.

→ Pooled estimate of $\pi$ When $\pi_1 = \pi_2 = \pi$ :

$$P = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2}$$

→ Test statistic for Hypotheses test about $\pi_1 - \pi_2$ :

$$Z = \frac{(P_1 - P_2)}{\sqrt{P(1-P)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

→ Reject $H_0$ :

→ If P-value $\leq \alpha$.

→ ① Lower tail test       ② upper tail test       ③ Two tail test

    $Z \leq -Z_\alpha$           $Z \geq Z_\alpha$           $|Z| \geq Z_{\frac{\alpha}{2}}$

## 11.1 Inferences about population variance :

$\delta^2$: Population variance

$\delta$ : Pop. st.dev

→ Sampling distribution of $\dfrac{(n-1)S^2}{\delta^2}$ has chi-squared distribution with $n-1$ degrees freedom.

$S^2$ : sample variance

$S$ : sample st.dev

→ Assuming :

1. The sample is Random.

2. The sample come from a Normal population.

→ $(1-\alpha)$ CI for $\delta^2 = \left( \dfrac{(n-1)S^2}{\chi^2_{\frac{\alpha}{2}}} \, , \, \dfrac{(n-1)S^2}{\chi^2_{1-\frac{\alpha}{2}}} \right)$ , where $df = n-1$ , $S^2 = \dfrac{\sum(x_i - \bar{x})^2}{n-1}$ .

→ Testing Hypothesis :

| ① lower tail test | ② upper tail test | ③ Two tail test. |
|---|---|---|
| $H_0: \delta^2 \geq \delta_0^2$ | $H_0: \delta^2 \leq \delta_0^2$ | $H_0: \delta^2 = \delta_0^2$ |
| $H_1: \delta^2 < \delta_0^2$ | $H_1: \delta^2 > \delta_0^2$ | $H_1: \delta^2 \neq \delta_0^2$ |

→ Test statistic :

$$\chi^2 = \dfrac{(n-1)S^2}{\delta_0^2} \quad , \quad df = n-1 \quad , \quad \delta_0^2 : \text{hypothesis value.}$$

→ Reject $H_0$ if :

→ P-value $\leq \alpha$ .

| ① lower tail test | ② upper tail test | ③ Two tail test. |
|---|---|---|
| $\chi^2 < \chi^2_{1-\alpha}$ | $\chi^2 > \chi^2_{\alpha}$ | $\chi^2 > \chi^2_{\frac{\alpha}{2}}$ or $\chi^2 < \chi^2_{1-\frac{\alpha}{2}}$ . |

## 11.2 Inferences about two population variances.

→ sampling distribution of $\frac{S_1^2}{S_2^2}$ when $\delta_1^2 = \delta_2^2$ has F distribution with $n_1 - 1$ df for the numerator and $n_2 - 1$ df for the denominator.

→ Assuming:

1. sample 1 and sample 2 random

2. sample 1 and sample 2 independent.

3. sample 1 and sample 2 are from Normal population.

4. $\delta_1^2 = \delta_2^2$.

→ Notations:

- $N_1$: size of pop. 1
- $N_2$: size of pop. 2
- $\delta_1^2$: variance of pop. 1
- $\delta_2^2$: variance of pop. 2
- $n_1$: size of sample 1
- $n_2$: size of sample 2
- $S_1^2$: variance of sample 1
- $S_2^2$: variance of sample 2

→ Testing Hypothesis:

① upper tail test:    lower $\delta_1^2$

$H_0: \delta_1^2 \leq \delta_2^2$

$H_1: \delta_1^2 > \delta_2^2$

② Two tail test:

$H_0: \delta_1^2 = \delta_2^2$

$H_1: \delta_1^2 \neq \delta_2^2$

→ Test statistic:

$$F = \frac{S_1^2}{S_2^2} \quad \text{with} \quad df_1 = n_1 - 1, \quad df_2 = n_2 - 1.$$

→ Reject $H_0$ if:

↝ p-value ≤ $\alpha$

↝ ① upper tail test

$F \geq F_\alpha$

② Two tail test.

$F \geq F_{\frac{\alpha}{2}}$

→ Note:

Population 1 is the population with higher sample variance.

# 12.1 Goodness of Fit and Independence :

→ $H_0 : \pi_1 = \pi_{10}$ , $\pi_2 = \pi_{20}$ , ..... , $\pi_K = \pi_{K0}$

$H_1 :$ The population proportion are not $\pi_1 = \pi_{10}, ... , \pi_K = \pi_{K0}$

→ $\pi_i :$ Population proportion of category $i$ .

$\pi_{i0} :$ Hypothesied value of the population of category $i$ . $i = 1, ... , K$.

→ $n :$ sample size

$f_i :$ observed frequencies.

$e_i :$ expeted frequencies .

→ $e_i = n \, \pi_{i0}$

→ $\sum\limits_{i=1}^{K} f_i = \sum\limits_{i=1}^{K} e_i = n$

→ Test statistic :

$$\chi^2 = \sum\limits_{i=1}^{K} \frac{(f_i - e_i)^2}{e_i}$$

→ Reject $H_0$ if P-value $\leq \alpha$ $\Big\}$ $df = K-1$  chi square table.

Reject $H_0$ if $\chi^2 \geq \chi^2_\alpha$ $\Big\}$ $e_i \geq 5$ $\forall i$ .

→ Remark

To use Goodness of fit test for multinomial populations we assume :

1. The sample taken is Random.

2. expected frequencies for all categories ; should satisfy the following $e_i \geq 5$ $\forall i$.

# 12.2: Test of Independence

→ Null and alternative hypotheses:

$H_0$: The Row variable and the column variable are independent.

$H_1$: The Row variable and the column variable are not independent.

→ We need to take a Random sample:

$f_{ij}$ : observed frequency

$e_{ij}$ : expected freq.

$n$ : # of Row

$m$ : # of columns

$$e_{ij} = \frac{(\text{Row } i \text{ table})(\text{column } j \text{ table})}{\text{sample size}}$$

→ Note: $\sum_j \sum_i f_{ij} = \sum_j \sum_i e_{ij} = \text{sample size}$.

→ Test statistic:

$$\chi^2 = \sum_j \sum_i \frac{(f_{ij} - e_{ij})^2}{e_{ij}} \qquad \text{with} \quad df = (n-1)(m-1).$$

Assuming: $e_{ij} \geq 5 \quad \forall i \; \forall j$.

→ Rejection Rule:

• Reject $H_0$ if $p$-value $\leq \alpha$.

• Reject $H_0$ if $\chi^2 \geq \chi^2_\alpha$.

→ Note: we use chi-square table

## 12.3: Goodness of fit test : Poisson and Normal distribution.

### Poisson distribution

→ $H_0$: The population has a Poisson dist.

$H_1$: The population doesn't have a Poisson dist.

→ Take a Random sample of size $n$

$f_i$ : observed frequencies $\quad \sum f_i = n$

$e_i$ : expected frequencies $\quad \sum e_i = n$

$$e_i = \frac{\mu^{x_i} \, e^{-\mu}}{x_i!} \cdot n \qquad \qquad \mu = \frac{\sum\limits^{K} x_i f_i}{\sum\limits_{i=1}^{K} f_i}$$

→ Test statistic:

$$\chi^2 = \sum_{i=1}^{K} \frac{(f_i - e_i)^2}{e_i} \qquad \text{with} \quad df = K-2 \qquad \qquad (\text{Assuming } e_i \geq 5 \; \forall i)$$

→ Rejection Rule:

• Reject $H_0$ if $p$-value $\leq \alpha$.

• Reject $H_0$ if $\chi^2 \geq \chi_\alpha^2$.

## Normal distribution

→ $H_0$: The population has a Normal distribution.

$H_1$: The population doesn't have a Normal distribution.

→ Take a Random sample of size $n$

$f_i$: observed frequencies.

$e_i$: expected frequencies.

→ Notation:

- $K$: # of categories.

$K = \dfrac{n}{5}$

- $e_i = 5 \quad \forall i$

→ Test statistic:

$$\chi^2 = \sum_{i=1}^{K} \frac{(f_i - e_i)^2}{e_i} \qquad \text{with} \quad df = K - 3.$$

→ Rejection Rule:

- Reject $H_0$ if P-value $\leq \alpha$.

- Reject $H_0$ if $\chi^2 \geq \chi^2_\alpha$.

## 13.2 : Analysis of variance : testing for the equality of K population means.

→ Testing for the equality of K pop. mean sample mean for treatment $j$ :

$$\overline{x}_j = \frac{\sum\limits_{i=1}^{n} x_{ij}}{n_j}$$

→ sample variance for treatment $j$ :

$$S_j^2 = \frac{\sum\limits_{i=1}^{n}(x_{ij} - \overline{x}_j)^2}{n_j - 1}$$

→ over sample mean :

$$\overline{\overline{x}} = \frac{\sum\limits_{j=1}^{n}\sum\limits_{i=1}^{n} x_{ij}}{n_T}$$

$$\overline{\overline{x}} = \frac{\sum\limits_{j=1}^{n} \overline{x}_j}{K} \qquad \text{if } n \text{ are equal}$$

> **Notations :**
>
> $x_{ij}$ : value of observation $i$ for treatment $j$.
>
> $n_j$ : number of observation for treatments $j$.
>
> $\overline{x}_j$ : sample mean for treatment $j$.
>
> $S_j^2$ : sample variance for treatment $j$.
>
> $S_j$ : sample standard deviation for treatment $j$.

→ Between treatments estimate of population variance :

. Mean square due to treatments .

$$MSTR = \frac{SSTR}{K-1} \qquad \text{where} \qquad SSTR = \sum\limits_{j=1}^{n} n_j (\overline{x}_j - \overline{\overline{x}})^2 .$$

→ Within treatments estimate of population variance :

. mean square due to error .

$$MSE = \frac{SSE}{n_T - K} \qquad \text{where} \qquad SSE = \sum\limits_{j=1}^{n} (n_j - 1) S_j^2 .$$

→ Test for the equality of K pop. means :

$H_0 : \mu_1 = \mu_2 = \ldots = \mu_K$

$H_1 :$ Not all population means are equal.

→ Test statistic :

$$F = \frac{MSTR}{MSE}$$

→ Rejection Rule :

p-value approach : Reject $H_0$ if $p$-value $\leq \alpha$

critical value approach : Reject $H_0$ if $F \geq F_\alpha$.

→ ANOVA table :

| Sorce of variance | df | SS | MS | F |
|---|---|---|---|---|
| Treatments | $k-1$ | SSTR | MSTR | $\dfrac{MSTR}{MSE}$ |
| Error | $n_T - k$ | SSE | MSE | |
| Total | $n_T - 1$ | SST | | |

# 13.3: multiple comparison procedures:

➡ **FLSD procedure:**

$$H_0^{ij}: \mu_i = \mu_j \qquad\qquad i \neq j$$
$$H_1^{ij}: \mu_i \neq \mu_j \qquad\qquad i,j \in \{1, \dots, k\}$$

➡ **Test statistic:**

$$\bar{X}_i - \bar{X}_j$$

➡ **Rejection Rule:**

Reject $H_0$ if $\quad |\bar{x}_i - \bar{x}_j| > LSD^{ij} \quad$ where $\quad LSD^{ij} = t_{\frac{\alpha_{CW}}{2}} \sqrt{MSE\left(\frac{1}{n_i} + \frac{1}{n_j}\right)}$

$$s.t \qquad \alpha_{CW} = \frac{\alpha_{EW}}{\binom{k}{2}} \qquad and \qquad df = n_T - k.$$

➡ **Bonferroni Adjustment:**

$\alpha_{EW}$: experiment wise Type I significance

$\alpha_{CW}$: comparision wise " "

- $\alpha_{EW} = \binom{k}{2} \alpha_{CW}$

- $\alpha_{CW} = \dfrac{\alpha_{EW}}{\binom{k}{2}}$

➡ **confedence interval.**

$(1-\alpha)$ CI for $\mu_i - \mu_j = (\bar{x}_i - \bar{x}_j) \mp LSD^{ij} \quad$ where $\quad LSD^{ij} = t_{\frac{\alpha}{2}} \sqrt{MSE\left(\frac{1}{n_i} + \frac{1}{n_j}\right)}$.

## 13.6 : Randomized Block design.

↳ contains Treatments and Blocks.

→ ANOVA table : Block Randomized Design :

| Source of variance | df | SS | MS | F |
|---|---|---|---|---|
| Treatments | $k-1$ | SSTR | $MSTR = \dfrac{SSTR}{k-1}$ | $F = \dfrac{MSTR}{MSE}$ . |
| Blocks | $b-1$ | SSBL | $MSBL = \dfrac{SSBL}{(b-1)}$ | |
| Error | $(k-1)(b-1)$ | SSE | $MSE = \dfrac{SSE}{(k-1)(b-1)}$ | |
| Total | $n_T - 1$ | SST | — | — |

* $b$ = # of Blocks .

* $BL$ = Blocks

* $n_T = kb$

* $SST = SSTR + SSBL + SSE$ .

→ Hypothesis :  $H_0 : \mu_1 = \mu_2 = \cdots = \mu_R$

$H_1 :$ Not all $\mu_j$ are equal .

→ Rejection Rule :

• critical value : reject $H_0$ if $F \geq F_\alpha$ with $df_1 = k-1$ , $df_2 = (k-1)(b-1)$.

• p-value : reject $H_0$ if p-value $\leq \alpha$ .

→ Notation and def :

- $\bar{x}_{i\cdot}$ : sample mean of Block $i$ , $i = 1, \ldots, b$ .

- $\bar{x}_{\cdot j}$ : sample mean of treatments $j$ , $j = 1, \ldots, K$ .

- $\bar{\bar{x}}$ : over all mean of all observation .

- $SST = \sum_{j=1}^{K} \sum_{i=1}^{b} (x_{ij} - \bar{\bar{x}})^2$ .

- $SSTR = b \sum_{j=1}^{K} (x_{\cdot j} - \bar{\bar{x}})^2$ .

- $SSBL = K \sum_{i=1}^{b} (x_{i\cdot} - \bar{\bar{x}})^2$ .

- $SSE = SST - SSTR - SSBL$ .

# 13.7: Factorial experiments

→ Notation

- $a$ = # of levels of factor A.
- $b$ = # of levels of factor B.
- $r$ = # of replications.
- $n_T$ = total number of observations taken in experiment, $n_T = abr$.

→ Hypothesis:

$H_0^A$ : means of factors A are equal

$H_1^A$ : means of factors A are not equal.


$H_0^B$ : means of factor B are equal.

$H_1^B$ : means of factor B are not equal.


$H_0^{AB}$ : Factor A and Factor B have no interaction.

$H_1^{AB}$ : Factor A and Factor B have an interaction.


→ Definitions

- $SST = SSA + SSB + SSAB + SSE = \sum\limits_{a=1}^{a} \sum\limits_{j=1}^{b} \sum\limits_{K=1}^{r} (x_{ijk} - \bar{\bar{x}})^2$.

- $SSA = br \sum\limits_{i=1}^{a} (\bar{x}_{i\cdot} - \bar{\bar{x}})^2$

- $SSB = ar \sum\limits_{j=1}^{b} (\bar{x}_{\cdot j} - \bar{\bar{x}})^2$

- $SSAB = r \sum\limits_{j=1}^{b} \sum\limits_{i=1}^{a} (\bar{x}_{ij} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{\bar{x}})^2$.

- $SSE = SST - SSA - SSB - SSAB$

→ ANOVA table : Two Factor Factorial experement.

| Source of variance | df | SS | MS | F | $F_\alpha$ | p-value |
|---|---|---|---|---|---|---|
| Factor A | $a-1$ | SSA | $MSA = \dfrac{SSA}{a-1}$ | $F^A = \dfrac{MSA}{MSE}$ | $F_\alpha$ with $a-1$, $ab(r-1)$ | |
| Factor B | $b-1$ | SSB | $MSB = \dfrac{SSB}{b-1}$ | $F^B = \dfrac{MSB}{MSE}$ | $F_\alpha$ with $b-1$, $ab(r-1)$ | |
| Interaction AB | $(a-1)(b-1)$ | SSAB | $MSAB = \dfrac{SSAB}{(a-1)(b-1)}$ | $F^{AB} = \dfrac{MSAB}{MSE}$ | $F_\alpha$ with $(a-1)(b-1)$, $ab(r-1)$ | |
| Error | $ab(r-1)$ | SSE | $MSE = \dfrac{SSE}{ab(r-1)}$ | — | — | — |
| Total | $n_T-1$ | SST | — | — | — | — |

$r \geqslant 2$

## 14.1 : simple linear Regression Model.

→ The simple linear Regression Model :

$$y = \beta_0 + \beta_1 X + \varepsilon$$

→ The simple linear Regression equation .

$$E(Y) = \beta_0 + \beta_1 X .$$

Note :
input
$X$ is a Variable
output
$y$ is a Random variable
$\varepsilon$ is a Random variable.

→ Estimated simple linear Regression equation .

$$\hat{y} = b_0 + b_1 X .$$

## 14.2  → Least square method :

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \qquad ; \text{ least square estimate for } \beta_1$$
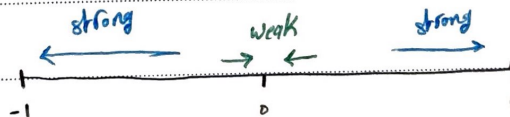
$$b_0 = \bar{y} - b_1 \bar{x} . \qquad ; \text{ least square estimate for } \beta_0 .$$

## 14.3 : coeffecient of determination .

→ correlation coeffecient :

$$\cdot \ r_{xy} = \frac{S_{xy}}{S_x S_y}$$

$$\cdot \ -1 \leq r_{xy} \leq 1$$

strong          weak          strong

←——————    →  ←    ——————→

-1                0                1

→ proposition :

$\cdot \ SST = SSR + SSE$            $\cdot \ SST = (n-1) s_y^2$

$\cdot \ SSR = b_1^2 (n-1) S_x^2$       $\cdot \ SSE = SST - SSR$ .

→ Coeffeacient of determination :

- $r^2 = \dfrac{SSR}{SST}$

- $0 \leq r^2 \leq 1$

- $r^2 = (r_{xy})^2 \rightarrow$ Coeffeacient of determination $= ($ correlation coeffecient $)^2$

- $r_{xy} = (\text{sign } b_1) \sqrt{r^2}$ .

## 14.5: Testing for significance

→ Model: $Y = \beta_0 + \beta_1 X + \varepsilon$.

→ $H_0: \beta_1 = 0$   means Not significance variable and Model.

$H_1: \beta_1 \neq 0$   means the Model and variable are significance.

→ Assuming: $E(\varepsilon) = 0$ , $Var(\varepsilon) = \delta^2$ , $\varepsilon$ independent , $\varepsilon$ Normal.

→ test statistic:

• t-test = $t = \dfrac{b_1}{S_{b_1}}$   with   $df = n-2$   } Two tail test.

Where $S_{b_1} = \sqrt{\dfrac{MSE}{(n-1) S_x^2}}$ .

اذا أعطانا $S_{b_1} = \dfrac{S}{\sqrt{(n-1) S_x^2}}$

→ Rejection Rule:

Reject $H_0$ if $|t| \geq t_{\frac{\alpha}{2}}$ .

→ Mean square Error (estimate of $\delta^2$)

$S^2 = MSE = \dfrac{SSE}{n-2}$ .

→ standard error of the estimate

$S = \sqrt{MSE} = \sqrt{\dfrac{SSE}{n-2}}$ .

→ Sampling distribution of $b_1$:

• $E(b_1) = \beta_1$

• $\delta_{b_1} = \dfrac{\delta}{\sqrt{\sum(x_i - \bar{x})^2}} = \dfrac{\delta}{\sqrt{(n-1) S_x^2}}$

• Distribution $b_1$ is Normal.

→ Estimated standared distribution of $b_1$ :

$$S_{b_1} = \frac{S}{\sqrt{(n-1)S_x^2}}.$$

→ $(1-\alpha)$ CI :

$$= b_1 \mp t_{\frac{\alpha}{2}} S_{b_1}.$$

2.qqd → F-test and ANOVA table.

| Source of variation | df | SS | MS | F |
|---|---|---|---|---|
| Regression | 1 | SSR | MSR | $\frac{MSR}{MSE}$ |
| Error | $n-2$ | SSE | MSE | |
| Total | $n-1$ | SST | — | |

↑
14.2

⤳ $MSR = SSR$

$$MSE = \frac{SSE}{n-2}$$

→ Rejection Rule :

- Reject $H_0$ if $F \geq F_\alpha$ with $df_1 = 1$ and $df_2 = n-2$.

- Reject $H_0$ if $p$-value $\leq \alpha$.

upper

14.6: using the estimated regression equation for estimation and predection

→ Point estimation : $\hat{y} = b_0 + b_1 X$.

→ Interval Estimation :

$$\hat{y}(P) = X_P$$

1. confidence interval for the mean value of $y$.

$$(1-\alpha) CI \text{ for } \underset{E(\hat{y}_p)}{y} = \hat{y}_P \mp S \, t_{\frac{\alpha}{2}} \sqrt{\frac{1}{n} + \frac{(x_P - \bar{x})^2}{\mathcal{E}(x_i - \bar{x})^2}} \quad \text{where } \hat{y}_p = b_0 + b_1 x_p \text{ and } df = n-2.$$

2. predection interval for $y$ :

$$(1-\alpha) PI \text{ for } y_P = \hat{y}_P \mp S \, t_{\frac{\alpha}{2}} \sqrt{1 + \frac{1}{n} + \frac{(x_P - \bar{x})^2}{\mathcal{E}(x_i - \bar{x})^2}} \quad \text{where } \hat{y}_p = b_0 + b_1 x_p \text{ and } df = n-2.$$

$$S = \sqrt{MSE}$$

14.7: computer solution:            Excel & spss

14.8:

Chapter 15 : Multiple Regression .

15.1 : Multiple Regression Model.

→ Model (Multiple linear Regression Model): $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$ .

→ Multiple linear Regression Equation : $E(Y) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$

→ Estimated multiple linear Regression : $\hat{y} = b_0 + b_1 X_1 + \dots + b_p X_p$

15.2 : least square Method :

$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \qquad\qquad \Big\} \quad By \; Excel.$$

15.3 : Multiple Coeffecient of determination :

→ $SST = SSR + SSE$ .

→ Multiple coeffecient of determinations : $R^2 = \dfrac{SSR}{SST}$ .

→ Adjusted Multiple coeffecient of determinations : $adj\,R^2 = 1 - (1 - R^2)\left(\dfrac{n-1}{n-p-1}\right)$

## 15.5 : Testing for significance.

→ Model : $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \varepsilon$.

→ $H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$

$H_1 :$ Not all $\beta_j$ are zero.

→ ANOVA Table :

| Source of variation | df | SS | MS | F |
|---|---|---|---|---|
| Regression | $p$ | SSR | MSR | $\dfrac{MSR}{MSE}$ |
| Error | $n-p-1$ | SSE | MSE | |
| Total | $n-1$ | SST | — | |

• $MSR = \dfrac{SSR}{p}$    • $MSE = \dfrac{MSE}{n-p-1}$    • $F = \dfrac{MSR}{MSE}$   with $df_1 = p$ and $df_2 = n-p-1$

• Reject $H_0$ if $F \geq F_\alpha$  or  p-value $\leq \alpha$

• F-test only one time.

---

→ $H_0 : \beta_j = 0$

$H_1 : \beta_j \neq 0$

→ Test statistic :   $t = \dfrac{b_j}{s_{b_j}}$

→ Reject $H_0$ if $|t| \geq t_{\frac{\alpha}{2}}$  or  p-value $\leq \alpha$    with $df = n-p-1$.

• t-test p-times.

→ $S = \sqrt{S^2} = \sqrt{MSE}$ : standared error of the estimate,

→ $\delta_{bj} = \sqrt{Var(bj)}$ : standard deviation of bj.

→ $S_{bj}$ = estimated standard deviation of bj.,

→ Molticolinearity :

input variable $x_1$, $x_2$, ... $x_p$

some times some $x_i$ is dependent on the other $x_j$'s, this case is known as multicolinearity.

→ Variance inflation factor (ViF)

$$ViF = \frac{1}{1 - R_j^2}$$

If $ViF(x_i) \geq 10$ then $x_i$ should be elimeneted.

• $R_j^2$ : Multiple coeffecient of determination for $x_i$ as a function of the other $x_j$'s.

• function means multiple regression.

Done ..