# Equation Summary

# Chapter 2:

① Measure of central tendancy:

1 → mean → Most accurate ~ affected by outliers

2 → median → middle accuracy ~ not affected by outliers

3 → mode → least accurate ~ affected by mode in far ends.

1 → mean : $\bar{x} = \dfrac{\Sigma x}{n}$        $\bar{x} = \dfrac{\Sigma fx}{\Sigma f}$

2 → Median = $Q_2$ = $n \times 50\%$ = $n^{th}$ term

Imp → $4^{th}$ term is $4^{th}$ term, $4.5^{th}$ term is $4^{th}$ and $5^{th}$ term.

3 → Mode = Most repeated number

∗ Coefficient of variation : $C.V = \dfrac{\sigma}{\mu}$

## ② Measure of spread:

1→ Range → max - min

2→ Inter-Quartile Range → $Q_3 - Q_1$

3→ Variance → $S^2$

4→ Standard deviation → $S$

$$\Rightarrow \quad S^2 = \frac{\sum x^2}{n-1} - \frac{\left(\sum x\right)^2}{n(n-1)}$$

## ③ Quantiles, percentiles and box-and-whisker:

$Q_1 = n \times 25\%$
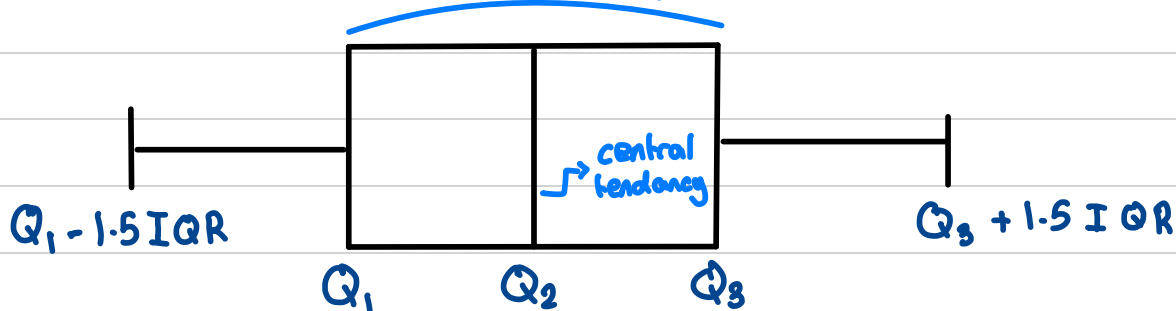
$Q_3 = n \times 50\%$

$P_k = (n) \times k\%$

If its a decimal, take next whole number. $7.15^{th} \to 8^{th}$

If whole, take that number and next one $10^{th} \to \frac{10^{th} + 11^{th}}{2}$

measure of spread



central tendency

$Q_1 - 1.5 IQR$

$Q_1$

n25%

↳ decimal → next
↳ whole → both

$Q_2$

n50%

↳ decimal → next
↳ whole → both

$Q_3$

n75%

↳ decimal → next
↳ whole → both

$Q_3 + 1.5 IQR$

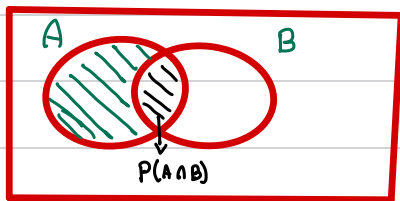Outliers:

Upper extreme $Q_3 + 1.5 IQR$

Lower extreme $Q_1 - 1.5 IQR$

# Chapter 3:

① $P(A \mid B) = \dfrac{P(A \cap B)}{P(B)}$    (conditional probability)

② $P(A \cap \bar{B}) = P(A) - P(A \cap B)$

basically, prob. of A and outside B
because $\bar{B} = 1 - B$.

So:



P(A∩B)

∴ $P(A \cap \bar{B})$
⇒ $\underset{\llcorner \text{green}}{P(A)} - \underset{\llcorner \text{black}}{P(A \cap B)}$

③ Total probability rule:

$P(A) = P(A \mid B) \times P(B) + P(A \mid \bar{B}) \times P(\bar{B})$

Given that $A \cap B$ & $A \cap \bar{B}$ are mutually exclusive (disjoint)

## ④ Baye's rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) \times P(A)}{P(B)}$$

$$= \frac{P(B|A) \times P(A)}{P(B|A) \times P(A) + P(B|\bar{A}) \times P(\bar{A})}$$
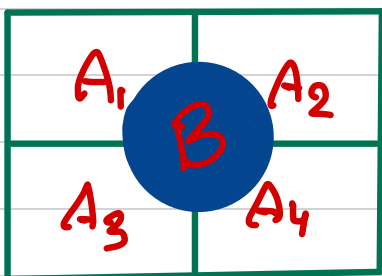
## ⑤ Total probability rule:

Given that events $A_1, \ldots, A_n$ are mutually exclusive

↳ can't occur at same time

↳ at least one event MUST occur

AND exhaustive

$$P(B) = \sum_{i=1}^{n} P(B|A_i) P(A_i)$$

$n = 4$

⑥ for any independent events:

$$P(A \cap B) = P(A) \times P(B)$$

⑦ for any two events:
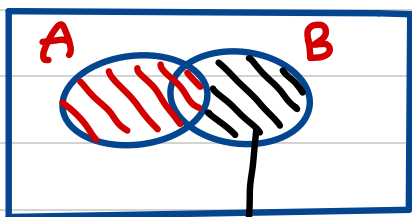
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

For mutually exclusive events:

$$P(A \cup B) = P(A) + P(B)$$
(because $P(A \cap B)$ does not exist in mutually exclusive events).

⑧ for two independent events:

$$P(A \cup B) = P(A) + P(B) \times P(\bar{A})$$



$P(B \cap \bar{A})$

$$P(A \cup B) = P(A) + P(B \wedge \bar{A})$$

$$= P(A) + P(B) \times P(\bar{A})$$

only unshaded region is
$P(A \cap B)$. $\therefore$ $P(\bar{A} \cup \bar{B}) = P(\overline{A \cap B}) = 1 - P(A \cap B)$

$P(A \cup B)$ is not shaded
$\therefore$ $P(\bar{A} \cap \bar{B}) = P(\overline{A \cup B}) = 1 - P(A \cup B)$

⑪ **Sensitivity** $= P(T^+ \mid D)$

⑫ **Specificity** $= P(T^- \mid \bar{D})$

⑬ $PV^+ = P(D \mid T^+)$

⑭ $PV^- = P(\bar{D} \mid T^-)$

⑮ $P(T^+ \mid D) = 1 - P(T^- \mid D)$

⑯ $P(T^- \mid \bar{D}) = 1 - P(T^+ \mid \bar{D})$

**Cheat code:**



Remember that first branch is $D / \bar{D}$ and second branch is $P/N$ where its probability is conditional.

# Chapter 4:

① Conditions for Binomial distribution:

→ A sample of independant trials n.

→ Only two possible outcomes, success p or failure q.

→ p and q are constant for each trial.

② $X \sim B(n, p) \Rightarrow P(X = a)$

$$\therefore P(X = a) = nC_a \times p^a \times q^{n-a}$$

$(q = 1 - p)$

③ Variance and standard deviation:

$$Var(x) = \sigma^2$$

$$Sd(x) = \sigma$$

$\Rightarrow$

$$\sigma^2 = npq$$

$$\sigma = \sqrt{npq}$$

④ Discrete Random variable:

| $x$ | 1 | 2 | 3 |
|-----|---|---|---|
| $P(X=x)$ | a | b | c |

$1 \rightarrow a + b + c = 1$

$2 \rightarrow E(x) = 1 \times a + 2 \times b + 3 \times c$

$3 \rightarrow Var(x) = 1^2 \times a + 2^2 \times b + 3^2 \times c$

# Chapter 5:

① $X \sim N(\mu, \sigma^2)$

$\Rightarrow \boxed{Z = \dfrac{X - \mu}{\sigma}}$  $\left\{ P(x > a) = P\left(Z > \dfrac{a - \mu}{\sigma}\right) \right.$

② $P(z > -a) = P(z < a)$

$P(z < -a) = P(z > a)$

$P(a < z < b) = P(z < b) - P(z < a)$

# Chapter 6:

Questions: They ask whats the probability of a sample mean $\overset{r> \bar{x}}{}$
or a sample proportion to be greater than $a$.
$\overset{L> \hat{p}}{}$

Or calculate confidence interval for sample mean or sample
proportion.

① **Standard error** of the mean:

$$Var(\bar{x}) = \frac{\sigma^2}{n}$$

$$\boxed{S.d(\bar{x}) = \frac{\sigma}{\sqrt{n}}}$$

$\frac{\sigma}{\sqrt{n}}$ is estimated by $\frac{S}{\sqrt{n}}$

**note:** $Var(\bar{x})$ & $sd(\bar{x})$ are <u>NOT</u> the same as $S^2$ & $S$, they are
not sample variance and s.d but they are s.d and variance
of the <u>set of sample means.</u>

② **Central limit theorem:** $\bar{x} \sim N(\mu, (\sigma/\sqrt{n})^2)$ $\overset{:}{\underset{)}{}}$ $P(\bar{x} > a)$
↓
$P(z > \frac{a - \mu}{\sigma/\sqrt{n}})$

$1 \to \sigma$ is known: $\boxed{z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}}$

$2 \to \sigma$ is unknown. <u>AND</u> $n < 30$ : $\boxed{t = \frac{\bar{x} - \mu}{S/\sqrt{n}}}$ degree of
Using S. freedom $= n-1$

$3 \rightarrow \sigma$ is unknown. AND $n > 30$ : $\boxed{Z = \dfrac{\bar{X} - \mu}{S/\sqrt{n}}}$

Using S.

③ **Confidence Interval** of the mean:

$$\mu = \bar{X} \pm E$$

E:

$1 \rightarrow \sigma$ is known : $\boxed{E = Z_{\frac{\alpha}{2}} * \dfrac{\sigma}{\sqrt{n}}}$

$2 \rightarrow \sigma$ is unknown AND $n \leqslant 30$ : $\boxed{E = t_{\frac{\alpha}{2}} * \dfrac{S}{\sqrt{n}}}$

Using S

$3 \rightarrow \sigma$ is unknown AND $n > 30$ : $\boxed{E = Z_{\frac{\alpha}{2}} * \dfrac{S}{\sqrt{n}}}$

Using S

**notes:**

↳ To find $Z_{\frac{\alpha}{2}}$. Use C.I % for D. Or $\frac{\alpha}{2}$ value for B.

↳ To find $t_{\frac{\alpha}{2}}$ : $1 - C.I\% \Rightarrow \alpha \Rightarrow \dfrac{\alpha}{2} \Rightarrow \dfrac{\alpha}{2} + C.I\% \Rightarrow t_{\frac{\alpha}{2}}$

Ex: 95% C.I $\Rightarrow$ 5% $\alpha$ $\Rightarrow$ 2.5% $\frac{\alpha}{2}$ $\Rightarrow$ 97.5% $\rightarrow$ look t-dist. $\rightarrow t_{\frac{\alpha}{2}} = 2.776$
      d.f = 4 (example)

↳ The bigger the C.I the more the error (smaller rejection region)

↳ The bigger the n, the less the error.

④ finding **sample size** $n$:

$$n = \left( Z_{\frac{\alpha}{2}} * \frac{\sigma}{E} \right)^2 \quad \rightarrow \text{ equation derived from } E.$$

↳ or $t_{\frac{\alpha}{2}}$ depending on conditions

⑤ Point estimation of $p$:

$$\hat{p} \sim N\left( p, \sqrt{\frac{pq}{n}}^2 \right)$$

$\hat{p} \Rightarrow$ sample proportion.

$p \Rightarrow$ population proportion.

$$\therefore \quad Z = \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}}$$

(always $Z$-dist)

⑥ Standard Error:

$$\text{S.E.} \sqrt{\frac{pq}{n}} \quad \text{Estimated by} \quad \sqrt{\frac{\hat{p}\,\hat{q}}{n}}$$

↳ S.E point estimate

⑦ Interval Estimate of $p$: (confidence interval)

→ point estimator

↳ for C.I, we use $\sqrt{\frac{\hat{p}\,\hat{q}}{n}}$. for C.I.I we use $\sqrt{\frac{pq}{n}}$

$$p = \hat{p} \pm E$$

$$E = Z_{\frac{\alpha}{2}} * \sqrt{\frac{\hat{p}\,\hat{q}}{n}}$$

# Chapter 7:

Questions: Testing for one sample mean or population

① Type 1 & Type 2 errors:

\* Type 1 error $\Rightarrow \alpha \Rightarrow$ Reject True $H_0$

**Mnemonic:**
$\alpha RT$
(art) $\alpha$ reject true $H_0$

\* Type 2 error $\Rightarrow \beta \Rightarrow$ Accept False $H_0$

$\beta AF$
(Baf) $\beta$ accept false $H_0$

② Testing of the hypothesis:

$\rightarrow$ numerical value

1$\rightarrow$ $H_0: \theta = \theta_0$ VS  \* $H_1: \theta > \theta_0 \Rightarrow$ Right tailed

\* $H_1: \theta < \theta_0 \Rightarrow$ Left tailed

\* $H_1: \theta \neq \theta_0 \Rightarrow$ Two tailed

2$\rightarrow$ Test Statistic

3$\rightarrow$ p-value method OR rejection region method

③ Test statistic for $\mu$:

1→ $\sigma$ is known: $\boxed{Z = \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}}}$

2→ $\sigma$ is unknown. AND $n < 30$ : $\boxed{t = \dfrac{\bar{X} - \mu}{S/\sqrt{n}}}$ degree of freedom $= n - 1$
Using S.

3→ $\sigma$ is unknown. AND $n > 30$ : $\boxed{Z = \dfrac{\bar{X} - \mu}{S/\sqrt{n}}}$
Using S.

④ Test statistic for $p$:

$$\therefore \quad Z = \dfrac{\hat{p} - p}{\sqrt{\dfrac{pq}{n}}}$$

$$\therefore \quad Z_{corr} = \dfrac{|\hat{p} - p| - \dfrac{1}{2n}}{\sqrt{\dfrac{pq}{n}}}$$

⑤ p-value method:

Use test stat to find p-value

* p-value > α      Accept $H_0$ & Reject $H_1$

* p-value < α      Reject $H_0$ & Accept $H_1$



| 0.001 | 0.01 | 0.05 ] p-values |

very highly    Highly   Significant        Significant        Not Significant
Significant

* Reject $H_0$ ⇒ Due statistically significant

⑥ Rejection region method:

Use α to find critical value. Use test statistic
to find z-value. If z-value in rejection region then reject $H_0$
else accept $H_0$.

# Chapter 8:

Questions: Two samples. Test hypothesis. for means only.

① Two Dependent Samples: (same n)

$1 \rightarrow$ d = after - before

$2 \rightarrow \bar{d} = \dfrac{\sum d}{n}$                    d.f = n - 1

$3 \rightarrow S^2_d = \dfrac{\sum d^2}{n-1} - \dfrac{(\sum d)^2}{n(n-1)}$

Test Statistic :

$$t = \dfrac{\bar{d}}{S_d / \sqrt{n}} \qquad n \leq 30$$

- - - - - - - - - - - - - - - - - - - - - - -

$$z = \dfrac{\bar{d}}{S_d / \sqrt{n}} \qquad n > 30$$

note that test stat is similiar to one sample mean, but d = 0 so we dont write it.

C.I $\Rightarrow$

$$d = \bar{d} \pm E \quad \Big| \quad E = t_{\frac{\alpha}{2}} * \dfrac{S_d}{\sqrt{n}}$$

② Two independent samples:

Test Statistic:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_p \cdot \sqrt{\frac{1}{n} + \frac{1}{m}}}$$

degree of freedom $= n + m - 2$

C.I:

$$\mu_1 - \mu_2 = (\bar{X}_1 - \bar{X}_2) \pm E$$

- - - - - - - - - - - - - - - - - - - - - -

$$E = t_{\frac{\alpha}{2}} * S_p * \sqrt{\frac{1}{n} + \frac{1}{m}}$$

Pooled Variance, $S_p \Rightarrow \sqrt{\frac{S_1^2(n-1) + S_2^2(m-1)}{n + m - 2}}$

# Chapter 10:

Questions: Testing a hypothesis for **two-sample** **proportion.**

$\rightarrow \hat{P}_1, \hat{P}_2$

① Normal Theory method, $z$ :

Conditions: \* $n\, p^*q^* > 5$   \& $m\, p^*q^* > 5$

\* Sample size is **large,** samples **discrete** and **independant**

① $H_0:$ $P_1 - P_2 = 0$   Vs   $H_1:$ $P_1 - P_2 \neq 0$

$\rightarrow$ two-tailed

(can also be right-tailed or left-tailed but $x^2$ is always right-tailed)

② **Test stat₁:**

$\hookrightarrow$ no correction

$$Z = \frac{(\hat{P}_1 - \hat{P}_2) - (P_1 - P_2)}{\sqrt{p^*\, q^*} \;*\; \sqrt{\frac{1}{n} + \frac{1}{m}}}$$

Is always zero.

$\hookrightarrow q^* = 1 - p^*$

**Test Stat₂:**

$\hookrightarrow$ with correction

$$Z_{corr} = \frac{\left|\hat{P}_1 - \hat{P}_2\right| - \left(\frac{1}{2n} + \frac{1}{2m}\right)}{\sqrt{p^*\, q^*} \;*\; \sqrt{\frac{1}{n} + \frac{1}{m}}}$$

pooled proportion, $\boxed{p^* = \dfrac{X + Y}{n + m}}$

$$\hat{P}_1 = \frac{x}{n} \qquad \hat{P}_2 = \frac{y}{m}$$

③ Either **p-value method** or **rejection region method**

② Contingency Table:

$$* \ \chi^2 = \sum \frac{(O-E)^2}{F}$$

$$* \ \chi^2_{corr} = \sum \frac{(|O-E| - 0.5)^2}{E}$$

note that we mainly use $\chi^2_{corr}$ for $2 \times 2$.

$$\chi^2 = \frac{(O_{11} - E_{11})^2}{E} + \frac{(O_{12} - E_{12})^2}{E} + \frac{(O_{21} - E_{21})^2}{E} + \frac{(O_{22} + E_{22})^2}{E}$$

## Observed (given):

|   | 1 | 2 |   |
|---|---|---|---|
| A | $x$ | $y$ | $x^{R_1} + y$ |
| B | $z$ | $u$ | $z^{R_2} + u$ |

$x + z$      $y + u$      $(x+y)+(z+u)$
$C_1$         $C_2$       or
                        $(x+z)+(y+u)$
                             $T$

## Expected

|   | 1 | 2 |
|---|---|---|
| A | $\dfrac{C_1 \times R_1}{T}$ | $\dfrac{C_2 \times R_1}{T}$ |
| B | $\dfrac{C_1 \times R_2}{T}$ | $\dfrac{C_2 \times R_2}{T}$ |

$$\therefore \quad E = \frac{R \times C}{T}$$

Degree of freedom: $(r-1) \times (c-1)$

↳ no of rows      ↳ no of columns

notes:

* $H_0$: independant variables. Two variables NOT associated

* $H_1$: dependent variables. Two variables ARE associated.

* for 2×2 table, all E must be greater than 5

* for R×C table,
Always the expected value is greater than 1.
For every 5 cells, only one value of E less
than 5 is allowed.

Ex:

| 5 | 9 | 121 |
|---|---|-----|
| 1 | 17 | 2 |

Doesn't
Satisfy
↵

| 5 | 9 | 121 |
|---|---|-----|
| 7 | 17 | 2 |

Satisfy
↵

* Always take right tailed test

# ③ Chi-squared goodness of fit test:

Probability × Observed total value = Expected value

1 → Continuity correction

2 → Find probability for cell

3 → multiply probability with $O_f$, get expected

4 → find expected for all cells.

5 → Chi-squared test. Accept or reject $H_0$.

↳ dont forget continuity correction

$$* \quad P(x < a) \times O_t$$

Degree of freedom for Chi-Squared goodniss -of- fit:

$$d.f = g - k - 1$$

g: number of groups/categories
↳ If E < 5 then for this category we join it with another category and it counts as one.

k: number of estimated parameters
↳ point estimates => $\bar{x}$, $s$, $\hat{p}$

# Chapter 11:

Questions: When they talk about ==Correlation== / ==Correlation coefficient==

① Covariance: (not very important)

$$Cov(x, y) = E((x-\mu_x)(y-\mu_y)) = E(xy) - \mu_x\mu_y$$

② ==Correlation coefficient==

population correlation coefficient $= \rho$

$$\rho = \frac{Cov(x,y)}{\sigma_x \; \sigma_y} \qquad \text{(equation not very important)}$$

Sample correlation coefficient $= r$

formula sheet

$$r = \frac{L_{xy}}{\sqrt{L_{xx} \; L_{yy}}} = \frac{S_{xy}}{S_x \cdot S_y}$$

$r$
‖ point estimator
↓
$\rho$

$L_{xx}: \; \sum x^2 - \frac{(\sum x)^2}{n}$

$L_{yy} = \sum y^2 - \frac{(\sum y)^2}{n}$

$L_{xy} = \sum xy - \frac{(\sum x)(\sum y)}{n}$

formula sheet

\* $-1 \leq r \leq 1$

→ closer to 1/-1 then stronger correlation

↳ 1 ⇒ +ve correlation
    -1 ⇒ -ve correlation

③ Statistical inference for hypothesis testing: Part 1

$H_0: \rho = 0$        Vs        $H_1: \rho \neq 0$

Test Stat = $\boxed{t = \dfrac{r\sqrt{n-2}}{\sqrt{1-r^2}}}$     In formula sheet.

④ Statistical inference for hypothesis testing: Part 2

$H_0: \rho = \rho_0$        Vs        $\rho \neq \rho_0$

      ↳ A non-zero number

Test stat:

1 → $z = \dfrac{1}{2}\ln\left(\dfrac{1+r}{1-r}\right)$     Formula sheet

2 → $z = \dfrac{1}{2}\ln\left(\dfrac{1+\rho}{1-\rho}\right)$     Get $\rho$ from statistical test

3 → $\lambda = (z - z_0)\sqrt{n-3}$     Formula sheet

note that the value of $\lambda$ is the value that you must compare with critical value on a z-distribution table.

⑤ Confidence Interval for population correlation coefficient and fisher's z-transformation :

① Find $(z_1, z_2)$
  ↳ fisher's

② Find $(\rho_1, \rho_2)$
  ↳ pop. corr. coeff.

① $z \pm E$

↳ $z = \dfrac{1}{2} \ln\left(\dfrac{1+r}{1-r}\right)$    ↳ $E = z_{\frac{\alpha}{2}} * \dfrac{1}{\sqrt{n-3}}$

∴ $z_{(1,2)} = z \pm z_{\frac{\alpha}{2}} * \dfrac{1}{\sqrt{n-3}}$

② $\rho = \dfrac{e^{2z}-1}{e^{2z}+1}$

* $\rho_1 = \dfrac{e^{2z_1}-1}{e^{2z_1}+1}$

C.I $= (\rho_1, \rho_2)$

* $\rho_2 = \dfrac{e^{2z_2}-1}{e^{2z_2}+1}$

# Chapter 12:

Questions : When they ask for testing of ==3 or more samples.==

① Hypothesis testing of ==More than 2 samples== using ==one way Anova Model.== :

$1 \rightarrow$ $H_0: \mu_1 = \mu_2 = \ldots \mu_k$     vs     $H_1: \mu_1 \neq \mu_2 \neq \ldots \mu_k$

  $\llcorner$ at least one pair not equal.

$2 \rightarrow$ Test Statistic $= \boxed{F = \dfrac{MS_B}{MS_W}}$

Formula sheet
$\Downarrow$

$*$ $\boxed{MS_B = \dfrac{SS_B}{k-1}}$     $\rightarrow$     $\boxed{SS_B = \sum n\bar{y}^2 - \dfrac{(\sum ny)^2}{N}}$

Degree of freedom $\Rightarrow$ ==k-1==

$\llcorner$ N is total Sample size

$*$ $\boxed{MS_W = \dfrac{SS_W}{N-k}}$     $\rightarrow$     $\boxed{SS_W = \sum (n-1)s^2}$

Degree of freedom $\Rightarrow$ ==N-k==

② ANOVA notes:

* D.F numerator ——— $= k-1$

$$\frac{MS_B}{MS_W}$$

* D.F denominator ——— $= N-k$

$$\boxed{SS_T = SS_B + SS_W}$$ ⇒ didn't see any questions in this.

③ Least Significant Difference test (LSD):

Lↄ I need some LSD after this exam 😵

Lↄ LSD is used to see which means are not equal to each other

Lↄ We compare two means individually, so we use the equation from chapter 8 for two independent samples.

Test Stat:

$$t = \frac{\bar{y}_1 - \bar{y}_2}{S.p \cdot \sqrt{\frac{1}{n} + \frac{1}{m}}}$$

$$S.p = \sqrt{MS_W}$$

↳ the $MS_W$ is used for all groups throughout so Sp is the same for any two means.

d.f $= N-k$

$$MS_W = \frac{SS_W}{N-k} = \frac{\sum(n-1)\cdot s^2}{N-k}$$

**Example:**

$H_0$: $\mu_1 = \mu_2$    Vs    $H_1$: $\mu_1 \neq \mu_2$

② Test stat $\Rightarrow$ $t = \dfrac{\bar{Y}_1 - \bar{Y}_2}{S.p \sqrt{\frac{1}{n} + \frac{1}{m}}}$    $Sp = \sqrt{MS_w}$

③ Use df and $\alpha$ to find critical values.

Thus either reject or accept $H_1$.

" لَا يُكَلِّفُ اللّٰهُ نَفْسًا إِلَّا وُسْعَهَا "