

Basic Statistics

3rd Semester

- Louisa Mandy Halim



Week 1

Three types of data

- quantitative, consists of numbers weights, heights
- Ranked, -!!— relative standing within a group
- qualitative, describing how they feel, look, wear

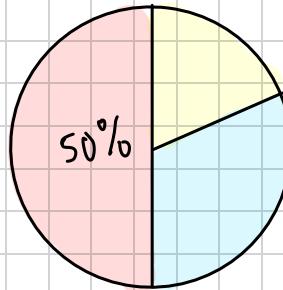
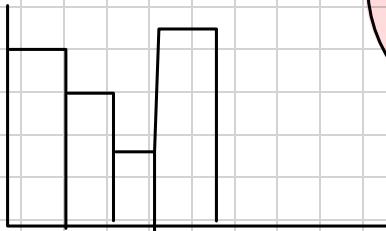
Level of measurement

- Sorting observations into different classes or categories
- classifying

independent & dependent variable

Week 2

- pie chart
- Bar chart



Comparing distributions

a bar chart of the number of people playing different card games on Wednesday and Thursday

Make simple and easy to understand designs for charts

example



Stem and leaf

22, 24, 26, 29
13, 14, 15, 16, 17, 19

6, 7

2		2, 4, 6, 9
1		3, 4, 5, 6, 7, 9
0		6, 7

leaf

stem

Histogram

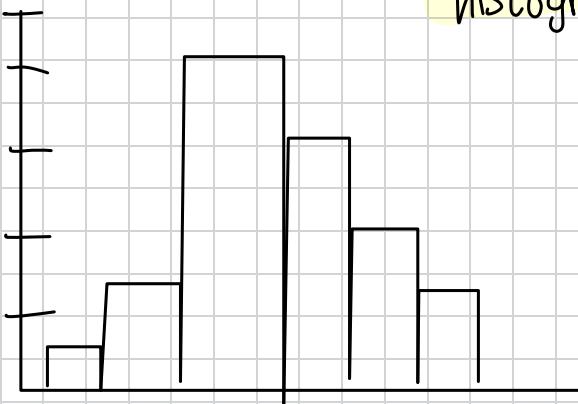
graphical method

displaying the shape of

distribution

- frequency table

histogram

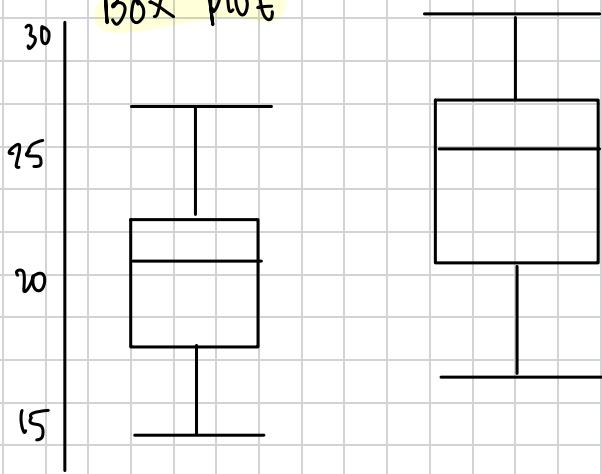


frequency polygons

lower	upper	Count

good choice to
display frequency
distribution

Box plot



- comparing genders
- identifying outliers and for comparing distributions

percentile

80% of the people are shorter than you

No cap

Exercise 01 — Sep 11, 2024

① Create a stem and leaf Display

Data set:

62, 65, 68, 70, 73, 75, 75, 78, 81, 83, 84, 85, 87, 89, 92, 95, 96, 98, 100

stem	leaf
6	2, 5, 8
7	0, 3, 5, 5, 8
8	1, 3, 4, 5, 7, 9
9	2, 5, 6, 8
10	0

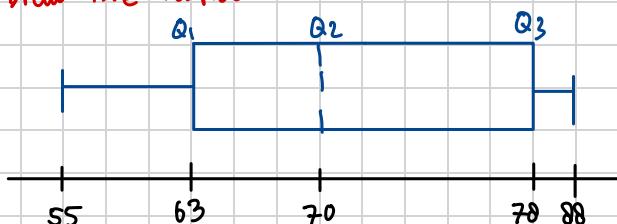
② Construct a box plot

Dataset: 55, 60, 62, 63, 65, 66, 68, 70, 72, 75, 77, 78, 80, 85, 88

a) Determine the five number summary

- minimum value = 55
- maximum value = 88
- $Q_1 = 63$
- $Q_2 = 70$ (middle)
- $Q_3 = 78$

b) Draw the boxplot



c) identify potential outliers

$$\begin{aligned} \text{IQR} &= Q_3 - Q_1 \\ &= 78 - 63 \\ &= 15 \end{aligned}$$

$$\begin{aligned} &\frac{1.5}{15} \times \\ &22.5 \end{aligned}$$

$$\begin{aligned} \text{upper fence} &= Q_3 + (1.5 \times \text{IQR}) \\ &= 78 + (1.5 \times 15) \\ &= 78 + 22.5 \\ &= 100.5 \end{aligned}$$

$$\begin{aligned} \text{lower fence} &= Q_1 - (1.5 \times \text{IQR}) \\ &= 63 - (1.5 \times 15) \\ &= 63 - 22.5 \\ &= 41.5 \end{aligned}$$

Since all data points lie within range 41.5 to 100.5. There are no outliers

Exercise 02 - Sep 25, 2024

① calculate the Trimean for a dataset below

Dataset :

10, 12, 15, 18, 21, 24, 27, 30, 33, 36, 39, 42, 45, 48, 50

Since the dataset is sorted already — find the median:

$$Q_1 = 18$$

$$Q_2 = 30$$

$$Q_3 = 42$$

Calculate the Trimean

$$\begin{aligned} \text{Trimean} &= \frac{Q_1 + 2 \times \text{Median} + Q_3}{4} \\ &= \frac{18 + 2(30) + 42}{4} \\ &= 120/4 = 30 \end{aligned}$$

The trimean
for this dataset
is 30

Trimean formula:

$$\frac{q_1 + (m \times 2) + q_3}{3}$$

② Geometric Mean

suppose that the population of a city changes over four years, with the following annual growth rates:

$$\begin{aligned} \text{Year 1: } &+ 5\% \rightarrow 1.05 \\ \text{Year 2: } &+ 10\% \rightarrow 1.1 \\ \text{Year 3: } &- 3\% \rightarrow 0.97 \\ \text{Year 4: } &+ 6\% \rightarrow 1.06 \end{aligned}$$

Geometric mean formula:

$$GM = \sqrt[n]{x_1 \cdot x_2 \cdot x_3 \dots}$$

Calculate the geometric mean of the growth rates to find the average population growth rate over these 4 years

$$GM = \sqrt[4]{1.05 \cdot 1.1 \cdot 0.97 \cdot 1.06}$$

$$= 1.0426$$

$$= 1.0426 - 1 = 0.0426 \times 100$$

$$= 4.26\% \rightarrow \text{the growth rate over 4 years}$$

(3)

Trimmed Mean

$$10\% \text{ of } 10 = \frac{10}{100} \times 10 = 1, \text{ remove 1 from each end}$$

10 ↗

~~65, 70, 72, 75, 80, 85, 90, 95, 100~~

calculate the 20% trimmed mean (Trim 10% from both sides)

$$\frac{70 + 72 + 75 + 80 + 85 + 90 + 95}{8} = \underline{\underline{82,375}}$$

add Conclusion

D More examples

a figure skating competition produces the following Scores:

6.0; 8.1; 8.3; 9.1; 9.9

$$\begin{aligned} \text{mean} &= \frac{6.0 + 8.1 + 8.3 + 9.1 + 9.9}{5} \\ &= \frac{41.4}{5} \rightarrow \underline{\underline{8.28}} \end{aligned}$$

a mean trimmed by 40% would equal 8.5 vs 8.28, which reduced the original mean by 0.22 points

Trim the mean by 40%

- remove the lowest 20% and the highest 20% of values, eliminating the scores: 6.0 & 9.9

$$\frac{8.1 + 8.3 + 9.1}{3} = \underline{\underline{8.5}}$$

Exercise 03 - Oct 02, 2024

- ① you have 8 people, and you need to select and arrange 4 of them in a row of photo. How many different ways can you arrange?

$$n=8 \quad r=4 \quad \left| \frac{n!}{(n-r)!} = \frac{8!}{4!} = \frac{40,320}{24} \rightarrow \underline{\underline{1,680 \text{ ways}}} \right.$$

- ② you have 7 books, and you want to choose 4 to take on a trip. How many different ways can you select the books?

$$n=7 \quad r=4 \quad \left| nCr = \frac{n!}{(n-r)! \cdot r!} \quad \frac{7!}{3! \times 4!} = \frac{7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1}{3! \times 4!} \right. \\ = \frac{210}{6} = \underline{\underline{35 \text{ ways}}}$$

- ③ a bag contains 10 red balls and 15 blue balls. If you randomly select 5 balls without replacement, what is the possibility that exactly 3 of the selected balls are red?

Combination formula : $nCr = \frac{n!}{(n-r)! \cdot r!}$

Hypergeometric formula : $p = \frac{k \binom{x}{k} \binom{N-x}{r-k}}{\binom{N}{r}}$

$$15C_2 = \frac{15!}{13! \cdot 2!} = \frac{15 \times 14}{2} \\ = 105$$

$$10C_3 = \frac{10!}{7! \times 3!} = \frac{10 \times 9 \times 8}{6}$$

$$25C_5 = \frac{25!}{20! \times 5!} = \frac{25 \times 24 \times 23 \times 22 \times 21}{120} \rightarrow 53,130$$

$$\text{Probability} \rightarrow \frac{10C_3 \times 15C_2}{25C_5} = \frac{120 \times 105}{53,130} = \underline{\underline{0.2372}}$$

More examples

in how many ways a committee consisting of 5 men and 3 women can be chosen from 9 men and 12 women

$$\rightarrow 9C_5 = \frac{9!}{4! \cdot 5!} \rightarrow \frac{9 \times 8 \times 7 \times 6}{4!} = \underline{\underline{126 \text{ ways}}}$$

$$\rightarrow 12C_3 = \frac{12!}{9! \cdot 3!} \rightarrow \frac{12 \times 11 \times 10}{6} = \underline{\underline{220 \text{ ways}}}$$

Total number of ways

$$\rightarrow 126 \times 220 = 27720$$

Exercise 04 - October 9, 2024

- ① find the percentage returns from an investment over 5 consecutive years, were:

Year 1: 10% → 0,1
 Year 2: 15% → 0,15
 Year 3: -5% → 0,95
 Year 4: 8% → 0,08
 Year 5: 12% → 0,12

$$GM = \sqrt[5]{0,1 \times 0,15 \times 0,95 \times 0,08 \times 0,12}$$

$$= 1,073$$

$$1,073 - 1 = 0,073 \times 100$$

$$= 7,3\%$$

- ② Box plot

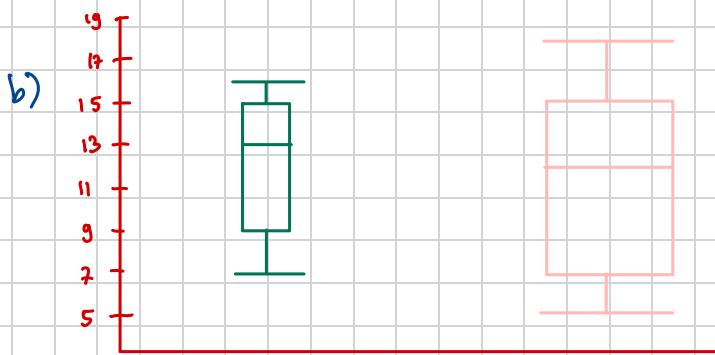
Dataset:

Group A: 7, 9, 12, 13, 14, 15, 16
 Group B: 5, 7, 8, 10, 12, 15, 18
 ↑ ↑ ↑

- a) calculate the five number summary

- min value = 7
- max value = 16
- $Q_1 = 9$
- $Q_2 = 13$
- $Q_3 = 15$

Group A



- c) compare the distributions of two groups based on the box plots.

- which group has a higher median? - Group A
- Are there any outliers?
no outliers

Group B



- ③ a card is drawn from a standard deck of 52 cards, and then a coin is flipped. What is the probability of drawing "king" from the deck and flipping a "Tail"?

1 Deck = 52 cards

Kings = 4

$$\frac{4}{52} \times \frac{1}{2} = \frac{1}{26} //$$

↑ ↓

four kings
out of the
total cards
tail side,
out of the
2 sides (coin)

- ④ Stem and leaf

Department X : 12, 14, 17, 19, 21, 24, 26, 28, 30, 32

Department Y : 13, 16, 18, 20, 23, 25, 27, 29, 31, 33

Back to back stem and leaf display

Department X		Department Y
2, 4, 7, 9		1 3, 6, 8
1, 4, 6, 8		2 0, 3, 5, 7, 9
0, 2		3 1, 3

- ⑤ The probability of getting exactly 3 heads when flipping a fair coin 5 times (where getting heads is considered a success)

$$5C_3 = \frac{5!}{2! \cdot 3!} = \frac{20}{2} = 10 \text{ ways}$$

$\left. \begin{array}{l} \frac{10}{32} = \frac{5}{16} \\ \text{Probability} \end{array} \right\}$

after 5 flips $\rightarrow 2^5 = 32$

Using the formula

$n = 5$ (number of trials)

$x = 3$ (number of successes)

$\pi = 0.5$ (probability of success)

$$\begin{aligned}
 & \frac{n!}{x!(n-x)!} \cdot \pi^x (1-\pi)^{n-x} \\
 &= 10 \times 0.5^3 \times 0.5^2 \\
 &= 10 \times \frac{5}{40} \times \frac{5}{10} \times \frac{5}{10} \times \frac{5}{10} \times \frac{5}{10} \\
 &= \frac{5}{16} \rightarrow \text{Probability}
 \end{aligned}$$

- (b) in a basketball game, a player has a free throw success rate of 80%. If the player takes 15 free throws, what is the probability that they make at least 12 successful free throws?

formula: $\binom{n}{x} \cdot p^x \cdot q^{n-x} \rightarrow \frac{n!}{(n-x)!x!} \cdot p^x \cdot q^{n-x}$

$$n = 15$$

$$x = 12, 13, 14, 15$$

$$p = 0,8$$

$$q = 0,2$$

o for $15 C_{12}$

$$\hookrightarrow \frac{15!}{3! \cdot 12!} \times 0,8^{12} \times 0,2^3 = 0,227$$

}

$$\text{Total} = 0,635$$

o for $15 C_{13}$

$$\hookrightarrow \frac{15!}{2! \cdot 13!} \times 0,8^{13} \times 0,2^2 = 0,236$$

}

the
Probability

o for $15 C_{14}$

$$\hookrightarrow \frac{15!}{1! \cdot 14!} \times 0,8^{14} \times 0,2^1 = 0,137$$

}

o for $15 C_{15}$

$$\hookrightarrow \frac{15!}{15!} \times 0,8^{15} \times 0,2^0 = 0,035$$

}

(7)

Hours of sunlight	Height (cm)	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$	$(y - \bar{y})^2$
2	10	-4	-10	40	16	100
4	15	-2	-5	10	4	25
6	20	0	0	0	0	0
8	25	2	5	10	4	25
10	30	4	10	40	16	100

$$\text{Total} = 30 \quad 100 \quad 0 \quad 0 \quad 100 \quad 40 \quad 250$$

$$\text{Mean} = 6 \quad 20 \quad 0 \quad 0$$

Calculate the Pearson correlation coefficient

$$\text{mean of } x = \downarrow 2, 4, 6, 8, 10$$

$$\text{mean of } y = \downarrow 10, 15, 20, 25, 30$$

$$\text{formula : } \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \cdot \sum (y - \bar{y})^2}} = \frac{100}{\sqrt{1000}} = \frac{100}{100} \rightsquigarrow \underline{\underline{1}}$$

Exercise 05 – October 23 Oct , 2024

- Find the standard deviation from those data

Scores : 70, 85, 75, 90, 88

$$\text{mean} \rightarrow \frac{70 + 85 + 75 + 90 + 88}{5} = 82,2$$

$$\frac{(70 - 82,5)^2 + (85 - 82,5)^2 + (75 - 82,5)^2 + (90 - 82,5)^2 + (88 - 82,5)^2}{5} = \frac{(-12,2)^2 + (2,8)^2 + (-4,2)^2 + (7,8)^2 + (5,8)^2}{5}$$

$$\begin{aligned} \text{variance} &= 53,76 \\ \text{standard deviation} &= \sqrt{53,76} \rightarrow \underline{\underline{7,33}} \end{aligned}$$

(2) Suppose a survey indicates that 30% of people prefer coffee over tea. If you randomly select 100 people, what is the probability that fewer than 25 people prefer coffee? Use z-table

$$n = 100$$

$$p = 30\% \rightarrow 0.3$$

$$q = 0.70$$

$$M = n \times p \\ = 100 \times \frac{3}{10} = 30$$

$$\sigma = \sqrt{100 \times 0.3 \times 0.7} \\ = \sqrt{21} \approx 4.58$$

because $25 \rightarrow 25 - 0.5 = 24.5$

$$z = \frac{24.5 - 30}{4.58} = \frac{-5.5}{4.58} \rightarrow -1.20$$

from z table, $-1.20 \rightarrow 0.1151$

The probability that fewer than 25 people prefer coffee is approximately 0.1151 or 11.51%

(3) You are conducting an experiment with 100 trials ($n=100$) and the probability of success in each trial is $p=0.4$, find the probability that at least 45 success will occur

$$n = 100$$

$$p = 0.4$$

$$q = 0.6$$

$$M = n \times p$$

$$= 100 \times 0.4$$

$$= 40$$

$$\sigma = \sqrt{100 \times 0.4 \times 0.6}$$

$$= 2\sqrt{6} \rightarrow 4.89$$

$$z = \frac{44.5 - 40}{4.89} \rightarrow \frac{4.5}{4.89} = 0.92$$

from z table $\rightarrow 0.8212$

$$1 - 0.8212 = 0.1788 \rightarrow 17.88\%$$

Exercise 06 — November 20, 2024

- ① a company claims their light bulb last 1000 hours on average.

950, 960, 970, 980, 1020, 1030, 990, 1010, 1000, 955 → mean = 990,5

Test whether the mean lifespan differs significantly from 1000 hours using $\alpha=0,05$

$$\text{t statistic formula : } t = \frac{m - \mu}{s/\sqrt{n}}$$

$s = 25,87$
 $\mu = 990,5$
 $m = 1000$
 $= \frac{990,5 - 1000}{25,87/\sqrt{10}} = -1,16$

Compare t statistic with Critical Value

Degrees of freedom : $n-1 \rightarrow 10-1=9$

at $\alpha=0,05$ (two tailed)

the critical t value is approximately $\pm 2,262$ (from t table)

Since $t = -1,16$ falls within the range $[-2,262; 2,262]$, we fail to reject the null hypothesis

- ② a fitness coach measures the weight of 8 clients before and after a 6 week training program.

Client	Before (kg)	After (kg)	Difference (d)
1	85	82	-3
2	78	75	-3
3	90	85	-5
4	76	74	-2
5	88	85	-3
6	81	78	-3
7	79	76	-3
8	92	89	-3

Conduct a paired t-test to determine if the training program significantly reduced weight. Use $\alpha=0,05$

Solution:

Null Hypothesis (H_0): The training program has no effect on weight ($\mu_d = 0$).
Alternative Hypothesis (H_1): The training program reduces weight ($\mu_d < 0$).

S_d or standard deviation formula:

$$S_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}} = 0,835$$

Conclusion: ($\alpha=0,05$)

Degrees of freedom = $8-1=7$

$t_{\text{crit}} = 2,365$

Since $-10,59 < -2,365$
reject null hypothesis
significant difference

$$\text{t statistic formula : } t = \frac{\bar{d}}{s_d / \sqrt{n}}$$

$\rightarrow -3,125 / 0,835 / \sqrt{8} = -10,59$

$$\text{Standard deviation} = \frac{\sum d}{n}$$

$\rightarrow -3 - 3 - 5 - 2 - 3 - 3 - 3 - 3 / 8 = -3,125$

3

A nutritionist wants to test if a new diet plan (Group A) significantly improves weight loss compared to a standard diet plan (Group B).

The following data was collected:

Group	Sample Size (n)	Mean Weight Loss (\bar{x})	Standard Deviation (s)
Group A (New)	25	8 kg	2
Group B (Standard)	25	6 kg	2.5

Perform an independent t-test to determine if the new diet plan significantly improves weight loss at a significant level of $\alpha = 0.05$

calculate t statistic

$$t = \frac{\bar{x}_A - \bar{x}_B}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}$$

$$t = \frac{8 - 6}{\sqrt{\frac{2^2}{25} + \frac{2.5^2}{25}}} \rightarrow \frac{2}{\sqrt{0.16 + 0.25}} \rightarrow \frac{2}{0.64} \rightarrow 3.13$$

$$\text{Degrees of freedom (df)} = n_A + n_B - 2 = 48$$

$$\text{from t table (one tail)} = 1.679$$

Since $t \text{ value} > \text{critical } t \text{ value}$
 $3.13 > 1.679$

reject null hypothesis

Exercise 07 – December 11, 2024

①

A researcher wants to compare the growth of plants under three types of fertilizers (A, B, and C).
The heights of the plants after 30 days (in cm) are:

Fertilizer A	Fertilizer B	Fertilizer C
15	20	25
16	22	27
14	19	26
15	21	28
17	20	24

mean 1 \bar{x}_A mean 2 \bar{x}_B mean 3 \bar{x}_C

$$n_f = 5 \\ I = 3$$

total mean

Does the type of fertilizer (A, B, or C) significantly affect plant growth (with $\alpha = 0.05$)?
Perform a one-way ANOVA to determine if fertilizer type affects plant growth.
Create a null hypothesis and alternative hypothesis first.

$$SSTR = 5 \left((15 - 20.6)^2 + (16 - 20.6)^2 + (14 - 20.6)^2 \right)$$

$$SSE = \text{each}$$

$$5 \left((15 - 15.4)^2 + (16 - 15.4)^2 + (14 - 15.4)^2 \right)$$

$$SSTR = \frac{(mean_1 - total\ mean)^2}{5} + \frac{(mean_2 - total\ mean)^2}{5} + \frac{(mean_3 - total\ mean)^2}{5}$$

Fertilizer A mean

$$\frac{15 + 16 + 14 + 15 + 17}{5} = \frac{77}{5} \rightarrow 15.4 \quad (\bar{x}_A)$$

Fertilizer B mean

$$\frac{20 + 22 + 19 + 21 + 20}{5} = \frac{102}{5} \rightarrow 20.4 \quad (\bar{x}_B)$$

Fertilizer C mean

$$\frac{25 + 27 + 26 + 28 + 24}{5} = 26 \quad (\bar{x}_C)$$

Overall means

$$\bar{x} = \frac{20 + 22 + 19 + 21 + 20 + 25 + 27 + 26 + 28 + 14}{15} = 20.6 \quad (\bar{x})$$

Sum of Squares

Sum of squares (SST)

$$SST = (15 - 20.6)^2 + (16 - 20.6)^2 + (14 - 20.6)^2 + (25 - 20.6)^2 + (27 - 20.6)^2 + (26 - 20.6)^2 + (28 - 20.6)^2 = 89.2$$

Sum of Squares between group (SSB)

$$SSB = 5 [(15.4 - 20.6)^2 + (20.4 - 20.6)^2 + (26 - 20.6)^2] \\ = 281.04$$

Sum of Squares within group (SSW)

in this case, $p < 0.001 < 0.05$
so we reject the null hypothesis

$$SSW = SS \text{ total} + SS \text{ between}$$

... the p value is extremely small ($p < 0.001$), which indicates very strong evidence against the null hypothesis. Therefore, we conclude that the type fertilizer has a significant effect on plant growth.

Conclusion

The F-static ($F = 24.35$) is significant effect on plant growth.
At least one fertilizer produces a different mean plant height.

②

A researcher wants to determine if there is an association between **plant type** and **fertilizer preference**. The researcher surveys 90 plants and records the following data:

Fertilizer	Plant Type A	Plant Type B	Plant Type C	Total
Fertilizer X	10	20	10	40
Fertilizer Y	15	10	5	30
Fertilizer Z	5	5	10	20
Total	30	35	25	90

Conduct a Chi-Square test of Independence whether plant type and fertilizer preference are independent at $\alpha = 0.05$.

calculate the expected Frequencies

$$\text{formula} = \frac{\text{Row total} \times \text{Column total}}{\text{Grand total}}$$

Fertilizer X

- plant type A: $\frac{30 \times 40}{90} = \frac{1200}{90} \rightarrow 13.33$

- plant type B: $\frac{35 \times 40}{90} = \frac{140}{9} \rightarrow 15.56 / 15.55$

- plant type C: $\frac{25 \times 40}{90} = \frac{100}{9} \rightarrow 11.11$

Fertilizer Z

- Plant type A:

$$\frac{30 \times 20}{90} = \frac{20}{3} = 6.66$$

- plant type B:

$$\frac{35 \times 20}{90} = \frac{70}{9} \rightarrow 7.77$$

Fertilizer Y

- plant type A: $\frac{30 \times 30}{90} = 10$

- Plant type B: $\frac{35 \times 30}{90} = \frac{35}{3} \rightarrow 11.66$

- plant type C: $\frac{25 \times 30}{90} = \frac{25}{3} \rightarrow 8.33$

- Plant type C:

$$\frac{25 \times 20}{90} = \frac{50}{9} \rightarrow 5.55$$

Fertilizer y:

- Plant type A:

$$\frac{(10 - 10)^2}{10} = 0$$

- plant type B:

$$\frac{(10 - 11.66)^2}{11.66} = 0.23$$

- plant type C:

$$\frac{(10 - 8.33)^2}{8.33} = 1.33$$

Fertilizer Z

- Plant type A

$$\frac{(5 - 6.66)^2}{6.66} = 0.41$$

- plant type B

$$\frac{(5 - 7.77)^2}{7.77} = 0.98$$

- Plant type C

$$\frac{(10 - 5.55)^2}{5.55} = 3.56$$

Fertilizer	Plant A	Plant B	Plant C	Total
X	13.33	15.56	11.11	40
Y	10	11.66	8.33	30
Z	6.66	7.77	5.55	20
Total	29.99 → 30	34.99 → 35	24.99 → 25	

Compute chi-Square statistic:

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

- plant type B

$$\frac{(20 - 15.56)^2}{15.56} = 1.26$$

- plant type C

$$\frac{(10 - 11.11)^2}{11.11} = 0.11$$

Summing all given values = 11.21

Degrees of freedom

$$df = (\text{number of rows} - 1) \times (\text{number of columns} - 1)$$

$$= (3-1) \times (3-1)$$

$$= 2 \times 2$$

$$df = 4$$

compare test statistics
to critical value

- if $\chi^2 \leq \chi^2_{\text{critical}}$, fail to reject H_0
- if $\chi^2 > \chi^2_{\text{critical}}$, reject H_0

Here:

$$\chi^2 = 11,21$$

$$\chi^2_{\text{critical}} = 9,488$$

$$\text{Since } 11,21 > 9,488$$

reject the
null hypothesis

(3) A professor wants to investigate whether the **type of programming language** (Python, Java, C++) and the **study method** (Self-Study, Instructor-Led) affects students' test scores. The professor records the test scores of students after completing a course under each combination of factors.

Language	Self-Study	Instructor-Led
Python	78, 82, 85	90, 88, 92
Java	72, 75, 74	85, 80, 84
C++	65, 68, 70	78, 75, 80

Perform a Two-Way ANOVA to determine if there are significant effects of programming language, study method, or their interaction on test scores.

Create all null hypotheses.

Use $\alpha = 0.05$

- find Grand / Group mean

$$\text{Python Self Study} = \frac{78 + 82 + 85}{3} = 81,66$$

$$\text{Python Instructor-Led} =$$

$$\frac{90 + 88 + 92}{3} = 90$$

$$\text{Java Self Study} = \frac{72 + 75 + 74}{3} = 73,66$$

$$\text{C++ Instructor} = \frac{78 + 75 + 80}{3} = 75,66$$

$$\text{C++ Self Study} = \frac{65 + 68 + 70}{3} = 67,66$$

$$\text{Java Instructor} = \frac{85 + 80 + 84}{3} = 83$$

$$\text{Grand mean } (\bar{x}) = \frac{81,66 + 73,66 + 67,66 + 90 + 83 + 77,66}{6} = 78,194$$

Group means

Python

$$(81,66 + 90) \div 2 = 85,83$$

Java

$$(73,66 + 83) \div 2 = 78,33$$

C++

$$(67,66 + 77,66) \div 2 = 72,66$$

Self

$$\text{total all} \div 9 = 74,33$$

Instructor

$$-11 = 83,55$$

Python & self study

$$81,66$$

python & instructor

$$90$$

Java and Self

$$73,66$$

Java & Instructor

$$83$$

C++ & self

$$67,66$$

C++ & Instructor

$$77,66$$

Compute Sum of Squares :

$$\text{Total} = \sum (x_{ij} - \bar{x})^2$$

$$\bar{x} = 78.94$$

- Sum of Squares for programming language (A)

$$SSA = 6 \cdot (85.83 - 78.94)^2 + 6 \cdot (78.33 - 78.94)^2 + 6 \cdot (72.66 - 78.94)^2$$

Python Java C++

$$SSA = 523,6956$$

- Sum of squares for factor study method (B)

$$SSB = 9 (74.33 - 78.94)^2 + 9 (83.55 - 78.94)^2$$

self study instructor

$$SSB = 382,5378$$

- Sum of Squares Within (error) Single numbers from tables - Group means

- SS Python & self study = $(78 - 81.66)^2 + (82 - 81.66)^2 + (85 - 81.66)^2 = 24,6668$

- SS Python & instructor-led = $(90 - 90)^2 + (88 - 90)^2 + (92 - 90)^2 = 8$

- SS Java & self study = $(72 - 73.66)^2 + (75 - 73.66)^2 + (74 - 73.66)^2 = 4,6668$

- SS Java & instructor-led = $(85 - 83)^2 + (80 - 83)^2 + (84 - 83)^2 = 14$

- SS C++ & self study = $(65 - 67.66)^2 + (68 - 67.66)^2 + (70 - 67.66)^2 = 12,9512$

- SS C++ & instructor-led = $(78 - 77.66)^2 + (75 - 77.66)^2 + (80 - 77.66)^2 = 12,6668$

$$SSE = 24,6668 + 8 + 4,6668 + 14 + 12,9512 + 12,6668$$

$$SSE = 76,9528$$

- Total Sum of Squares all single numbers from table - \bar{x}

- SS Total = $(78 - 78.94)^2 + (82 - 78.94)^2 + (85 - 78.94)^2 + (72 - 78.94)^2 + (75 - 78.94)^2 + (74 - 78.94)^2 + (65 - 78.94)^2 + (68 - 78.94)^2 + (70 - 78.94)^2 + (90 - 78.94)^2 + (88 - 78.94)^2 + (92 - 78.94)^2 + (85 - 78.94)^2 + (80 - 78.94)^2 + (84 - 78.94)^2 + (78 - 78.94)^2 + (75 - 78.94)^2 + (80 - 78.94)^2 = 984,9444$

$$\begin{aligned}
 SS_{\text{interaction}} &= SS_{\text{Total}} - SSA - SSB - SSE \\
 &= 984,9444 - 523,6956 - 382,5378 - 76,9528 \\
 &= 1,7582
 \end{aligned}$$

Degrees of freedom

$$\begin{aligned}
 df_A &= 2 \\
 df_B &= 1 \\
 df_{\text{interaction}} &= 2 \\
 df_{\text{within}} &= 12 \\
 df_{\text{total}} &= 17
 \end{aligned}$$

Mean Squared and F-Statistics

Mean Squares and F-Statistics:

$$\begin{aligned}
 MS_A &= SS/df = 523.4394/2 = 261.7197 \\
 MS_B &= 191.3611 \\
 MS_{A+B} &= 1.0852 \\
 MS_E &= 38.3332 \\
 F_A &= MS_A / MS_E = 261.7197 / 38.3332 = 40.9652 \\
 F_B &= 59.9045
 \end{aligned}$$

$$F_{A+B} = 0.1656$$

Decision:

p-value Programming Language for
 $F = 40.965$, $df = (2, 12)$ at $\alpha = 0.05$ is 0.00000435
 p-value Study Method for
 $F = 59.9045$, $df = (1, 12)$ at $\alpha = 0.05$ is 0.00000527
 p-value Interaction for
 $F = 0.1656$, $df = (2, 12)$ at $\alpha = 0.05$ is 0.84928886

Conclusion:

Significant main effects of programming language on test scores. (p-value < α)
 Significant main effects of study method on test scores. p-value < α
 No Significant interaction between language and study methods. p-value > α

Source of Variation	SS	df	MS	F	P-value
Sample (study)	523,444444	2	261,722222	40,9652174	4,3476E-06
Columns (program)	382,722222	1	382,722222	59,9043478	5,26602E-06
Interaction	2,11111111	2	1,055555556	0,16521739	0,849605144
Within	76,66666667	12	6,388888889		
Total	984,944444	17			

Week 4

22

Probability of a single event

$$\text{Probability} = \frac{\text{Possible outcomes}}{\text{total outcomes}}$$

- probability to shuffle an ace in a one deck card
- the chance of dices
- probability of two (or more) independent events
- flipping a coin

Exercise 1:

Let's consider a scenario where you're trying to predict the exam score based on the number of hours studied, but the sample data has more variability due to other factors like individual learning styles, distractions, or even exam difficulty.

Hours of Study (X)	Exam Score (Y)
1	52
2	59
3	62
4	64
5	72
6	80
7	74
8	83
9	91
10	89

1. Find the regression equation ($Y' = bX + A$).
2. After that, predict the exam score for someone who studied for 7 hours. Compare the results with the data in the table, explaining why the predictions are not the same.
3. Then, predict the exam score for someone who studied for 11 hours.

Exercise 2:

You are given the following sample data showing the height (in inches) and weight (in pounds) of 5 people:

Height (X)	Weight (Y)
60	140
62	145
64	160
66	170
68	155

1. Find the regression equation ($Y' = bX + A$).
2. After that, predict the weight for someone who is 70 inches tall.

Exercise 1

Louisa Mandy Halim - 2702325552

①

Hours of Study (X)	Exam Score (Y)
1	52
2	59
3	62
4	64
5	72
6	80
7	74
8	83
9	91
10	89

X · Y	x^2	y^2
52	1	2.704
118	4	3.481
186	9	3.844
256	16	4.096
360	25	5.184
480	36	6.400
518	49	5.476
664	64	6.889
819	81	8.281
890	100	7.921

$$\Sigma = 55$$

$$\Sigma = 726$$

$$\Sigma = 4343 \quad \Sigma = 385 \quad \Sigma = 54.276$$

$$b_0 = \frac{((\Sigma y)(\Sigma x^2) - (\Sigma x)(\Sigma xy))}{(n(\Sigma x^2) - (\Sigma x)^2)}$$

$$b_0 = \frac{(726 \times 385) - (55 \times 4343)}{(10 \times 55) - (55)^2}$$

$$b_0 = \frac{279.510 - 238.865}{825}$$

$$b_0 = \frac{40.645}{825} \Rightarrow 49.3$$

$$b_1 = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2}$$

$$b_1 = \frac{(10 \times 4343) - (55 \times 726)}{(10 \times 385) - (55)^2}$$

$$b_1 = \frac{43430 - 39930}{3850 - 3025} = 4.24$$

$$y' = b_0 + b_1 x$$

$$y' = b_1 x + A$$

$$y' = \underline{\underline{4.24 x + 49.3}}$$

② to predict the student who studied for 7 hours

$$y' = 4.24 x + 49.3$$

$$x = 7 \rightarrow 7 \text{ hours}$$

$$y' = 4.24(7) + 49.3$$

$$y' = 78.98$$

The regression equation captures the general trend and minimizes errors, but may not precisely align with each individual data point.

③ Predict after studying for 11 hours

$$y' = 4.24 x + 49.3$$

$$x = 11 \rightarrow 11 \text{ hours}$$

$$y' = 4.24(11) + 49.3$$

$$y' = 95.94$$

Exercise 2

①

x	y	x^2	$x \cdot y$
60	140	3600	8400
62	145	3844	8990
64	160	4096	10240
66	170	4356	11220
68	155	4624	10540
$\Sigma =$	320	770	20520
			49390

$$b_0 = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2}$$

$$b_0 = \frac{5(49.390) - (320)(770)}{5(20.520) - (320)^2}$$

$$b_0 = \frac{550}{200} = 2,75$$

$$b_1 = \frac{(\sum y \cdot \sum x^2) - (\sum x \cdot \sum xy)}{n(\sum x^2) - (\sum x)^2}$$

$$b_1 = \frac{(770 \times 20.520) - (320 \times 49.390)}{200}$$

$$b_1 = -22$$

$$y' = b_0 x + A$$

$$y' = b_0 + b_1 x$$

$$y' = 2,75x - 22$$

② $y' = 2,75(70) - 22$

$y' = 170,5$

