

RM

— / —

UNIT 3

* Probability Sampling

* Stratified Sampling

* Disproportionate Sampling : When strata differ in size and variability, it is reasonable to take larger samples from variable strata

$$\frac{n_1}{N_1 \sigma_1} = \frac{n_2}{N_2 \sigma_2} = \dots$$

$$n_i = \frac{n \cdot N_i \sigma_i}{N_1 \sigma_1 + N_2 \sigma_2 + \dots + N_k \sigma_k}$$

σ → standard deviations

N → sizes

n → sample size (total)

n_i → sample size (of i strata)

$N_1 = 5000$ $N_2 = 2000$ $N_3 = 3000$
 $\sigma_1 = 15$ $\sigma_2 = 18$ $\sigma_3 = 5$
 $n = 84$ $n_1, n_2, n_3 = ?$

$$n_i = \frac{n \cdot N_i \sigma_i}{N_1 \sigma_1 + \dots + N_k \sigma_k}$$

i) $n_1 = \frac{n \cdot N_1 \sigma_1}{N_1 \sigma_1 + N_2 \sigma_2 + N_3 \sigma_3} = \frac{84(5000 \times 15)}{5k15 + 2k18 + 3k5}$

$$n_1 = \frac{6300000}{126000} = 50$$

ii) $n_2 = \frac{n \cdot N_2 \sigma_2}{126000} = \frac{84(2000 \times 15)}{126000} = 24$

iii) $n_3 = \frac{n \cdot N_3 \sigma_3}{126000} = \frac{84(3000 \times 5)}{126000} = 10$

IF stratum has cost C

$$n_i = \frac{n \cdot N_i \sigma_i}{\sqrt{\frac{N_1 \sigma_1}{C_1} + \dots + \frac{N_k \sigma_k}{C_k}}}$$

* Mean, Standard Deviation, Variance

X	f	fx	$[x - \bar{x}]^2$	$f[x - \bar{x}]^2$
2	40	80	4	160
4	30	120	0	0
6	20	120	4	80
8	10	80	16	160
Σx	100	400		400

* MEAN = $\bar{x} = \frac{\sum fx}{\sum f} = \frac{400}{100} = 4$

* STANDARD DEVIATION = $s = \sqrt{\frac{\sum f[x - \bar{x}]^2}{\sum f}} = \sqrt{\frac{400}{100}} = 2$

$$s = 2$$

* VARIANCE = $s^2 = 4$

* Coefficient of Variance = $\frac{s}{\bar{x}} = \frac{2}{4} = 0.5$

Standard Deviation = SD = s = σ

FORMULA

* Confidence Interval for the Mean

$$\mu = \bar{x} \pm z \frac{\sigma}{\sqrt{n}}$$

$$\mu \text{ lies between } \left[\bar{x} - z \frac{\sigma}{\sqrt{n}}, \bar{x} + z \frac{\sigma}{\sqrt{n}} \right]$$

* Confidence Interval for the Average

$$\mu = \bar{x} \pm z \frac{\sigma}{\sqrt{n}} \quad \begin{cases} N-n \\ N-1 \end{cases}$$

Mean = $\bar{x} = 10.455$

Sample = $n = 55$

S.D = $\sigma = 7.7$

Construct 90% confidence interval (CI)

Given for 90%, $z = 1.645$

$$\bar{x} - z \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + z \frac{\sigma}{\sqrt{n}}$$

$$10.455 - (1.64) \frac{7.7}{\sqrt{55}}, \quad 10.455 + (1.64) \frac{7.7}{\sqrt{55}}$$

$$[8.545, 12.365]$$

DATA ANALYSIS

*

Problem Definition

Well defined problem is 50%
of the project done

*

Plan selection

*

Hypothesis

*

Data Analysis

(hypothesis rejected or accepted)

Goal: to gain info from data

#

EDA

- Descriptive statistics

- Data visualization

-

Data Cleaning

-

Correlation Analysis

-

Dimensionality Reduction

-

Data Transformation

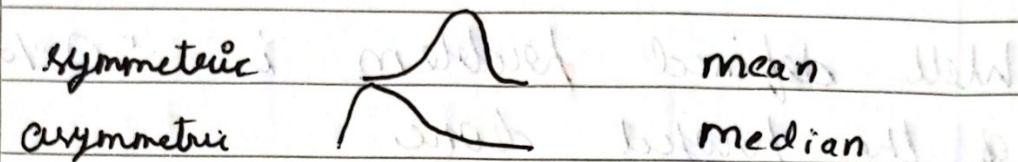
(normalization, age 1-80 to 0-1 scale)

-

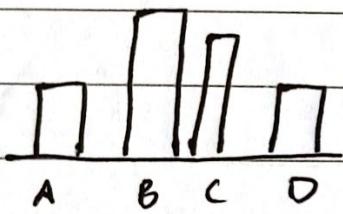
Summary Statistics

* Data Distribution

[Mean, Median, Mode] (Central Tendency)



Mode : which value occurs most frequently



B max times

$$2 \times \text{Mean} + \text{Mode} = 3 \times \text{Median}$$

- Quartile

Q1 : 25% percentile

Q2 : 50% percentile

Q3 : 75% percentile

IQR : Q3 - Q1

- [SD] (Variation)

$$s = \sqrt{\frac{(x - \bar{x})^2 + \dots + (y - \bar{y})^2}{n - 1}}$$

$$\sigma = s^2$$

~~Unit~~ 3 → ~~3~~ question

Unit 3 → ² ~~3~~ questi

Unit 4 → 3 question

* Statistical Estimation

- Confidence Intervals

Size (n) ↑, Confidence ↑

Problems 1, 2, 3, 4, 5

- T - Distribution

- Hypothesis Testing

Using t-distribution

When $n < 30$

And $\sigma = \text{Unknown}$

i) Data = {3, 1, 3, 7, 4, 2, ...}

$$n = 14$$
$$C = 99\%$$

ii) Calculate mean (\bar{x})

iii) Calculate SD (σ)

$$t_{\left(\frac{\alpha}{2}, n-1\right)} \Rightarrow \frac{\alpha}{2} = \frac{1-0.99}{2} = 0.005$$
$$n-1 = 14 - 1 = 13$$

iv) Confidence Interval, (μ)

$$\bar{x} - t_{0.005} \frac{\sigma}{\sqrt{n}}, \bar{x} + t_{0.005} \frac{\sigma}{\sqrt{n}}$$

Data = 3, 1, 3, 2, 5, 1, 2, 1, 4, 2, 1, 3, 1, 1
n = 14
C = 99%

a) MEAN, $\bar{x} = \frac{30}{14} = 2.143$

b) SD, $\sigma = \sqrt{\frac{\sum(\bar{x}-x)^2}{n-1}} = \sqrt{12.7} = 13$

$$\sigma = \sqrt{0.977} = 0.988$$

c) $t_{\frac{\alpha}{2}, n-1} = \frac{1-0.99}{2} = 0.005$
 $n-1 = 13$

$$t_{(0.005, 13)} = 3.012$$

d) $\bar{x} - 3.012 \frac{0.988}{\sqrt{14}} = 1.348$

$$\bar{x} + 3.012 \frac{0.988}{\sqrt{14}} = 2.94$$

* Population Proportion

Sample = $n = 200$
Proportion = $x = 80$
 $C = 95\%$

$$\rightarrow P = \frac{x}{n} = \frac{80}{200} = 0.4$$

$$\rightarrow Z = 95\% = 1.96$$

$$\rightarrow \mu = P + Z \sqrt{\frac{P(1-P)}{n}}$$

$$= 0.4 \pm 1.96 \sqrt{\frac{0.4(1-0.4)}{200}}$$

$$= [0.332, 0.468]$$

HYPOTHESIS TESTING

* NULL Hypothesis $\Rightarrow \mu = 20$

* ALTERNATE Hypothesis $\Rightarrow \mu \neq 20$

* Significance level = α = threshold

If $p \leq \alpha$, reject null hypothesis

If $p > \alpha$, accept null hypothesis

I. PARAMETRIC TESTS

1) t - test

1) One Sample

Compare mean of group to mean of population

$\rightarrow H_0 \Rightarrow \mu = \mu_0$ ($\frac{\text{sample mean}}{\text{mean}} = \frac{\text{population mean}}{\text{mean}}$)

$\rightarrow H_A \Rightarrow \mu \neq \mu_0$

$$\rightarrow t = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \rightarrow df = n - 1$$

\rightarrow find p value

$\rightarrow p \leq \alpha$ (n) $p \geq \alpha$

Data = {10.5, 9, 6.5, 8, 11, 7, 7.5, 8.5, 9.5, 12}

$$n = 10$$

$$\bar{x} = 12$$

→ Mean of sample, $\bar{x} = 8.95$

Mean of population, $\mu = 12$

$$\rightarrow H_0 \Rightarrow \mu = \bar{x}$$

$$H_A \Rightarrow \mu \neq \bar{x}$$

$$\rightarrow t = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{8.95 - 12}{1.802 / \sqrt{10}} = -5.35$$

$$\sigma = \sqrt{\frac{\sum (\bar{x} - \mu)^2}{n-1}} = \sqrt{\frac{29.225}{9}} = \sqrt{3.247} = 1.802$$

$$\rightarrow df = 9 \text{ and } \text{Significance} = 5\% = 0.05$$

$$\rightarrow \text{Tabulated Value } (9, 0.05) = 2.2622$$

\Rightarrow Calculated value (t) within $[-2.2, +2.2]$

$\therefore t < -2.2$, Reject NULL hypothesis

CONCL There is significant difference in mean of sample and population

One Sample, Two Occassion

BP before : 180 200 - - -

BP after : 140 145 - - -

Difference : 40 55 80 85

$$\rightarrow \sum d = 465 \text{ and } n = 8$$

$$\bar{d} = \frac{465}{8} = 58.125$$

$$\rightarrow \sigma = \sqrt{307.49} = 17.54$$

$$\rightarrow t = \frac{\bar{d}}{\sigma / \sqrt{n}} = 9.37$$

\rightarrow Tabulated Value (0.05 and 7)

___ / ___ / ___

* $Z - \text{test}$

- When $n > 30$ (Sample is large)

$$- Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

$\mu = 40$

$$n = 50 \quad \bar{x} = 45 \quad \sigma = 20$$

$$\alpha = 0.02$$

$$\rightarrow H_0 : \mu \leq \underline{\underline{40}}$$

$$H_A : \mu \neq \underline{\underline{40}}$$

$$\rightarrow Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{45 - 40}{20 / \sqrt{50}} = 1.767$$

$$\rightarrow \alpha = 0.02, \text{ df} = 49 \Rightarrow Z_{0.02} = 2.05$$

$$\rightarrow p > z$$

Do not reject H_0

— / — / —

#

$$\mu = 166.17$$

$$\sigma = 5.89$$

$$n = 144 \quad \bar{x} = 164.65$$

$$\alpha = 0.05$$

→ H_0 : Slum area residents are shorter ($\mu > 166$)

H_A : Slum area residents are ~~not~~ short

$$\rightarrow Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{164.65 - 166.17}{5.89/\sqrt{144}} = -3.097$$

$$\rightarrow Z_{0.05} = [-1.65, 1.65]$$

⇒ Reject null hypothesis

CHI SQUARE TEST (χ^2)

#		\square	\blacksquare	Total
	Male	30	20	
	Female	25	25	50
Total		55	45	100

→ H_0 : Gender and Vote are independent
 H_A : Variables are dependant

→ Contingency

$$E_{\text{Male}, \square} = \frac{55 \times 50}{100} = 27.5 \quad E_{F, A} = 27.5$$

$$E_{\text{Male}, \blacksquare} = \frac{45 \times 50}{100} = 22.5 \quad E_{F, B} = 22.5$$

→ Chi Square

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

$$\chi^2 = \frac{(30 - 27.5)^2}{27.5} + \frac{(20 - 22.5)^2}{22.5} + \dots + \frac{(25 - 22.5)^2}{22.5}$$

$$\chi^2 = 1.01$$

$$\rightarrow df = (2-1)(2-1) = 1$$

$$\rightarrow \alpha = 0.05 \text{ and } df = 1 \Rightarrow 3.841$$

$$\rightarrow \chi^2 < 3.841 \Rightarrow \text{Do not reject } H_0$$

Is Die Throw unbiased

Num	1	2	3	4	5	6
Frequency	16	20	25	14	29	28

$$\rightarrow E = 132 \times \frac{1}{6} = 22$$

$$\rightarrow \chi^2 = \sum \frac{(O - E)^2}{E} = \frac{(16-22)^2}{22} + \dots$$

$$\chi^2 = 9$$

$$\rightarrow \alpha = 0.05$$

$$df = 6-1 = 5$$

$$\rightarrow P = [-11, 11]$$

$\rightarrow \therefore$ Accept H_0

\therefore Die is unbiased