5min

Probability and probability Distributions

Statistic :← Central tende

① Descriptive

② Inferential

⟶ Hypothesis Testing.

Notes

population mean = SM (Sample) is mostly
Incorrect.

$$\underline{PM = SM} \quad \cancel{X}$$

$$\boxed{PM = SM \pm \text{Margin of Error}} \quad ✓ \quad \underline{\text{Can be Correct.}}$$

5 years

* Always try to keep Nu. of Sample 2000

at least 200 or

≥ 200

* According to Central limit theorem each
Samples must have 30 data points.

* Continuous data      col
                        |
                numerical value
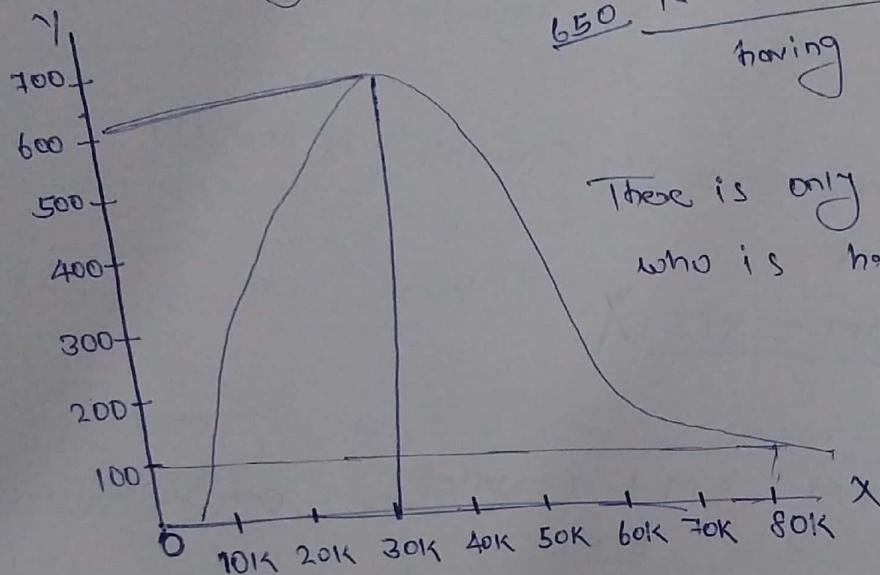                    float
                   decimal

Kde — kernal density estimate.

↳ when you have very extreme value at Right. for your graph is positively skew.

1:16

(C)   POSITIVELY skewd



650 No of employee having only 30k Salary

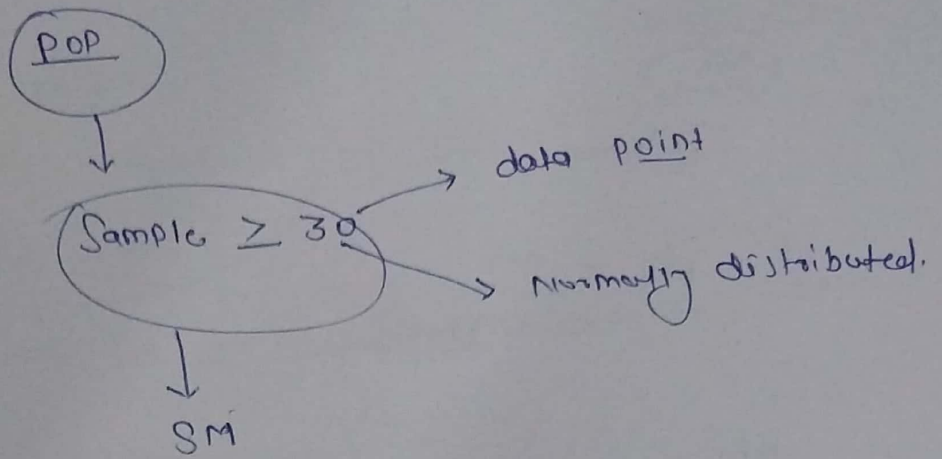There is only 100 employee who is having 80 Salary.

So, Analysis of graph tell there very let's people who is getting test high Salary & maximum People are getting Less Salary.

So, the graph is negatively skewd.

* opposite of this graph is positively skewd and either Normal / no skwed

* The bigger the Sample is, the better it is for us.

And

POP

↓

Sample ≥ 30 → data point

→ Normally distributed.

↓

SM

* if we don't have much resources in this case
Sample Size more than or equal to 30. is
also greate

$$SS = \geq 30$$

[with Central limit → proven
theorem. ← thragh.

Ram Sanyog Kumaa
19,1 14, 25, 15, 7

→ PM = SM ± Margin of Error

↓

This is the things so we are
studying the Central Limit
theorem.   1:30 hour: Min

Ram

35
1  0
2  0
3  0

Average income
of Indram families

PM = SM ± Margin of Error

Mean

PM = 25/1

= 25-1 to 25+1

= 24 to 26 → Range / Interval

And

POP

↓

Sample ≥ 30 → data point

→ Normally distributed.

↓

SM

* if we don't have much resources in this case
Sample size more than or equal to 30. is
also greater

$$SS = ≥ 30$$

[With Central limit → Proven
theorem. ← through.

Ram Sanyog kumaa
19,19,25, 15,7

→ PM = SM ± Margin of Error

Ram

This is the things so we are
studying the Central Limit
theorem.    1:30 hour : Min

Average income ————————→ ⟨35⟩  ①  O

of Indian families                              ②  O

PM = SM ± Margin of Error              ③  O

Mean

PM = 25/1

= 25-1 to 25+1

= 24 to 26 → Range / Intensay

(PM)

fish

(xived) POP

If    PM = SM

Hook

$$\frac{PM}{fish} \longrightarrow \frac{SM \pm Margin}{Net}$$

90%  $\longrightarrow$  25K $\pm$ 1K $\longrightarrow$ 24K to 26K

95%  $\longrightarrow$  25K $\pm$ 1.5K $\longrightarrow$ 23.5 to 26.5K

99%  $\longrightarrow$  25K $\pm$ 2.5K

confidence Interval

90% + $\frac{24K \ to \ 26K}{\downarrow \ Interval}$

CI = $\boxed{SM \pm Margin}$

PM

Formulas to find out confidence Interval

$$CI = \bar{x} \pm Z^* \times \frac{S}{\sqrt{n}}$$
         $\downarrow$
        (It)

will tell

Margin of error

$\overline{x}$ : Sample Mean

S : Standard deviation of the Sample.

n : Sample Size

$z^*$ : Z-score for a Certain Confidence level.

| $z^*$ | Confidence level | |
|-------|------------------|---|
| 1.65  | 90 %             | ✓ |
| 1.96  | 95%              | ✓ |
| 2.58  | 99 %             | ✓ |

## MVV I

### Business Problem

Estimate whether the mean lead content in maggi packets is within the allowed range or not?

Allowed range $=$ 2.5 PPM ( part per million )

Lets take value

n = 100 ( maggi packet)

$\overline{x}$ = 2.3 PPM

S = 0.3 PPM

① 2.20 PPM

② 2.43 PPM

③ 2.37 PPM

④ 2.21 PM

⋮

⑩⑩ 2.28 PPM

→ 50 M

→ 100 M

↓

PM

↓

X

SM ────→ =2.23 PM

$$CI = \bar{x} \pm z^* \frac{S}{\sqrt{n}}$$

$$= 2.3 \pm 2.58 \times \frac{0.3}{\sqrt{100}}$$

$$= 2.3 \pm 2.58 \times 0.03$$

$$= 2.3 \pm 0.07$$

$$= 2.3 + 0.07 \text{ to } 2.3 - 0.07$$

POPM

$$CI = \boxed{2.23 \, pm \text{ to } 2.37 \, Ppm.}$$

$$\underline{99\% \quad \text{Confidence}}$$

$$z^*$$
$$\boxed{\begin{array}{l} 90\% \\ 95\% \\ 99\% \end{array}}$$

POP

maggi

(100M)

But of this

→ 1

For Study :- Statquest

Financial Analyst

Stock (Apple) → $\boxed{50-60\%}$

$$\underline{90\% \quad \text{confidence}}$$

→ $\underline{45-75\%}$
$$95\%$$

→ $\boxed{30-80\%}$
$$99\%$$

Hypothesis Testing = 16

→

So    Inferential Statistic is

way to        PM ——→ SM ± Margin
                ↑ ←        ←

CI = SM ± Margin

= 29.5 ± 1

= 28.5 to 30.5

(the) average Life = 30month
                    → is - Correct

SM

↑

500

SS

① 28

② 27.2

③ 31

④ 30.5

⋮

500 / SM = 29.5



68.27%

95.45%

99.73%

$\mu - 3\sigma$    $\mu - 2\sigma$    $\mu$    $\mu + \sigma$    $\mu + 2\sigma$    $\mu + 3\sigma$

65% 

55k   60k                              70k   75k

Salary
20K People

$M = 65K$

$SD = 5K$.

2ey

$\dfrac{65K - 5K}{M - 1 \times SD}$
$\downarrow$
60K

$M + 1 \times SD$
$65 + 5K$
$\underline{\underline{70K}}$

→ 68.27% People earning
60K — 70K Salary.

* Notes
Through Central limit theorem when my $SS \geq 30$
and we plot these 30 values then I will
normally distributed.

50 Types
$\overrightarrow{SS}$ → POM POM 30
$SM = \underline{\underline{31.8}}$ ] So population mean is
correct

* when my rejection level is going to be One
Side.
We called it one tailed test.

* When my rejection level is going to be so
both the side we called it
two tailed test.

\* Rejection area should be of 5% of intese area.

$$Z = \frac{x - \mu}{\sigma}$$

$Z =$ Standard Score

$x =$ Observed value

$\sigma =$ Standard deviation of sample

Null hyp. Cont: $=$, $<=$, $>=$

Alternate Hyp. Co: $\neq$, $>$, $<$

according to Question

Null hyp $<=$

Alt $\geq$

Note :- Our aim is always to reject the null hypothesis.

Step 1 - B PS

Step 2 NH & AH
try to the reject the null hypothesis through some test.

Step: 3 Sample data collection

Step 4 SM $\longrightarrow$ Z - value of SM
$$\left( Z = \frac{x - \mu}{\sigma (SD)} \right)$$

⑤    SM    VS    PM    (with margin)

⑥    if    SM's   z-   value   fall into   redection region. (population mean claim is rejected)
       or Null hypotheis will be rejected

---

     Rejection

     Region

(old)   Method

Steps ⌐ (Critical value)

⟹ SM ⟹ SM z value ⟹ Critical values, critical value VS SM z value

⟹ then deciding we are rejecting the null hypotheis or not.

＊ Notes   Significante Leave gives as 5% rejection area.

＊ Notes :- when our p-value is less than the Significant value means %.5 or 0.05 then only we can reject the null hypothesis.

or   Let's

      z- test

        SD

        ↓

     p- value (0.03)

      3%

＊ Always Go with p-value

     $\dfrac{\bar{x}-\mu}{\dfrac{\sigma}{\sqrt{}}}$

    97% → A NH

Hae it mean 97% proof is against the null hypo-
thesis and 3% in favour of null hythesis.

**And**

we can reject the null hypothesis when we have
the proof against the null hypothesis, and
it should be _major [%]_ percentage.

→ Journey of hypothesis testing:-

1.) BPS

② NH & AH

③ Sample data ⟶ Hyp test.

④ Sample data ⟶ Hypothesis Test.

⑤ proof against NH in the form of p-value
or Critical value.

⑥ we should use the p-value to make the
final decision.

⑦ if the p-value is $\leq \underset{(\alpha)}{SL}$ (5%)

alpha
α

β → Beta        Then the proof in favour of NH
is too weak & the proof against NH is
is too strong. So we can reject the null hypo-
—thesis.

$$P-value \nless (0.05)$$

We fail to reject NH

0.43     $\nless$ 0.05    so we don't have enough

In the         proof to reject the null hypo-

favour of NH          thesis.

## Z- test (Assumptions)

① Sample size $\geq$ 80

② Data should be normally distributed

③ Population Standard Deviation must be given to us.

④ Sample should be randomly collected.

(Note) for 2

for checking this the data is normally distri-
buted or not ploting will not help.

so to verify this we will use another

test which

$\longrightarrow$ Shapiro wilks test.

T—test

① Sample size $\geq$ 30 or Sample size $\leq$ 30.
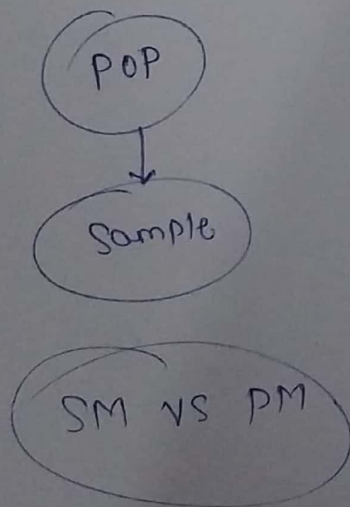
→ so data is normally distributed (Asumption)

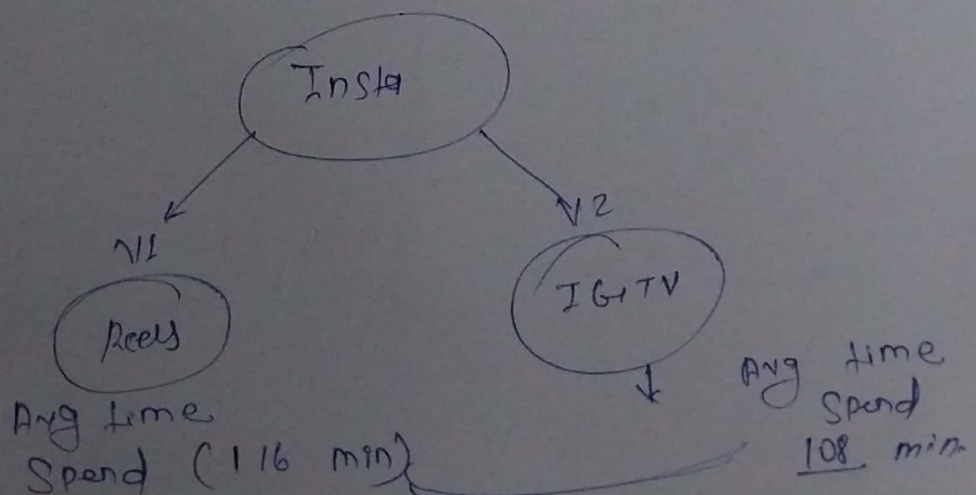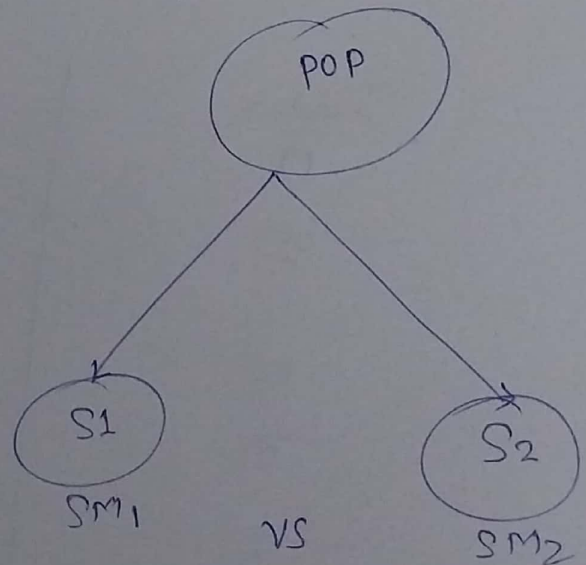② Data should be normally distributed.

↳ Shapiro will test for normality.

③ Sample Standard deviation.

NOTE) T—test is most use test rather than any other.

One Sample Test or Z—test

Two Sample t—test or Z—test



POP

↓

Sample

SM vs PM

POP

S1
SM1

S2
SM2

VS

Insta

√1

Reels

Avg time
Spend (116 min)

√2

IG+TV

↓ Avg time
Spend
108 min

ANOVA $\longrightarrow$ Numerical Data

for 3 or more than 3 Samples.

$\longrightarrow$ Analysis of variances

## Chi - Squared Test of Ind.

$\downarrow$

use on Categorical Data

Suppose

Starbucks

| Gender | Preference |
|--------|------------|
| M | T |
| F | C |
| F | C |
| M | C |
| M | T |
| M | C |
| F | T |
| F | C |
| M | T |

Example of ANOVA

Pharma



M1                M2                M3

40                40                40

People

Book fo DS in Shikha Sir, Grithub.

Statistical thinking for programmers.