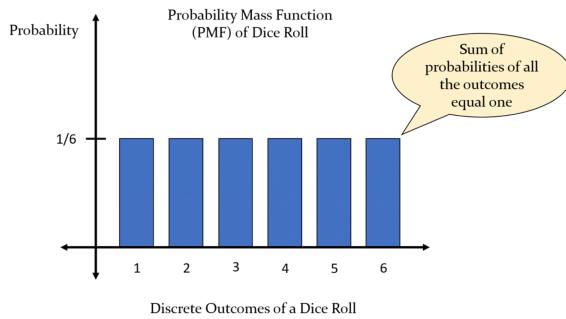


Probability Distribution →

A mathematical function that gives the probabilities of occurrence of different possible outcomes of an experiment. In simple words, it lists out all possible outcomes in the sample space with their probabilities.

Ex → Throwing a dice →

1	2	3	4	5	6
$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Discrete Probability Distribution →

Model the probabilities of random variables that can have discrete values as outcomes.

Ex → Let's take random variable X to be the sum of two die throws - X can take values $(2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12)$.

$$P(2) : (1,1) \rightarrow 1/36$$

$$P(3) : (1,2), (2,1) \rightarrow 2/36$$

$$P(4) : (2,2), (1,3), (3,1) \rightarrow 3/36$$

$$P(5) \rightarrow 4/36$$

$$P(6) \rightarrow 5/36$$

$$P(7) \rightarrow 6/36$$

:

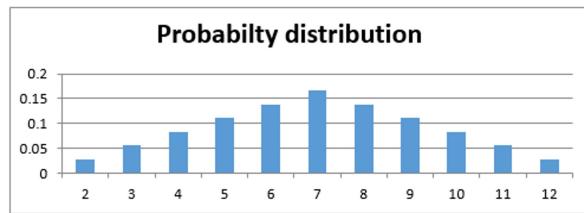
$$P(11) \rightarrow 2/36$$

$$P(12) \rightarrow 1/36$$

$$P(\text{success}) = 5/36$$

$$P(\text{Failure}) = 31/36$$

The probability function for a discrete random variable is the 'probability mass function'. It shows the exact probabilities for a particular value of the random variable.



Frequently used Discrete probability distribution in Data Science:

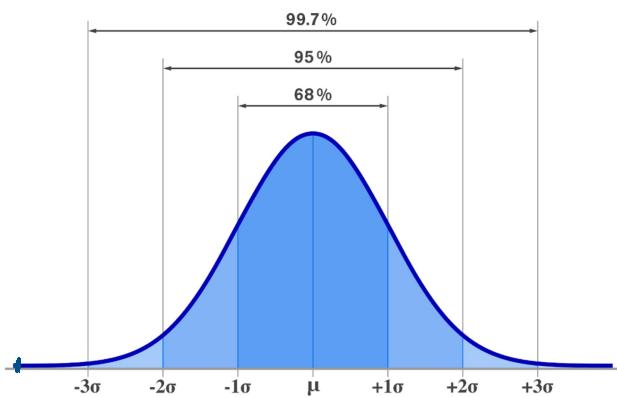
- Bernoulli Distribution.
- Binomial Distribution.
- Poisson Distribution.
- Multinomial Distribution.

Continuous Probability Distribution →

Describes the probabilities of the possible values of a continuous random variable.

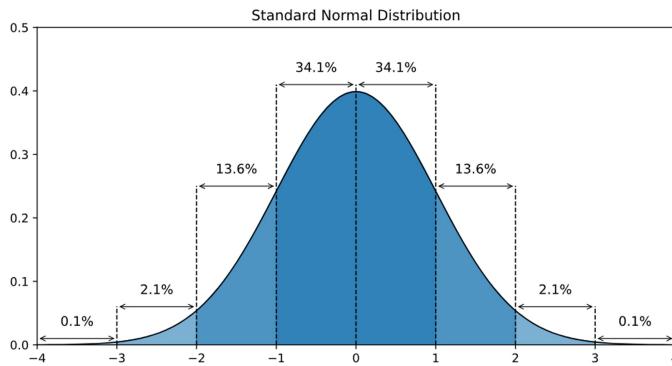
Normal Distribution →

One of the most used distribution characterized by its symmetric, bell-shaped curve.

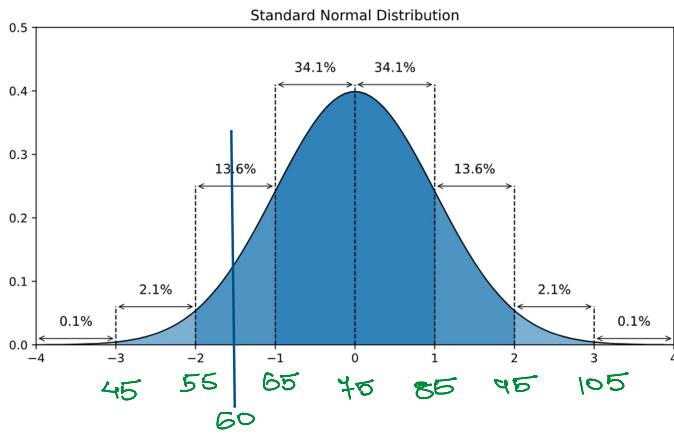


Many natural phenomenon also show normal distribution.

Standard Normal Distribution & Z-scores →



Normal Distribution can be represented as Standard Normal Distribution so that we can calculate the percentage of data present within a range.



Suppose we have a data on salary of employees of a company where mean = ₹5k and standard deviation = 10k.

Here we want to find how many employees are earning upto 60k.

We will use Z-score for the same:

$$Z = \frac{x - \text{Mean}}{\text{SD}}$$

$$= \frac{60 - 50}{10}$$

$$= -1.5$$

Now we check for this value in Z-score

$$= -1.5$$

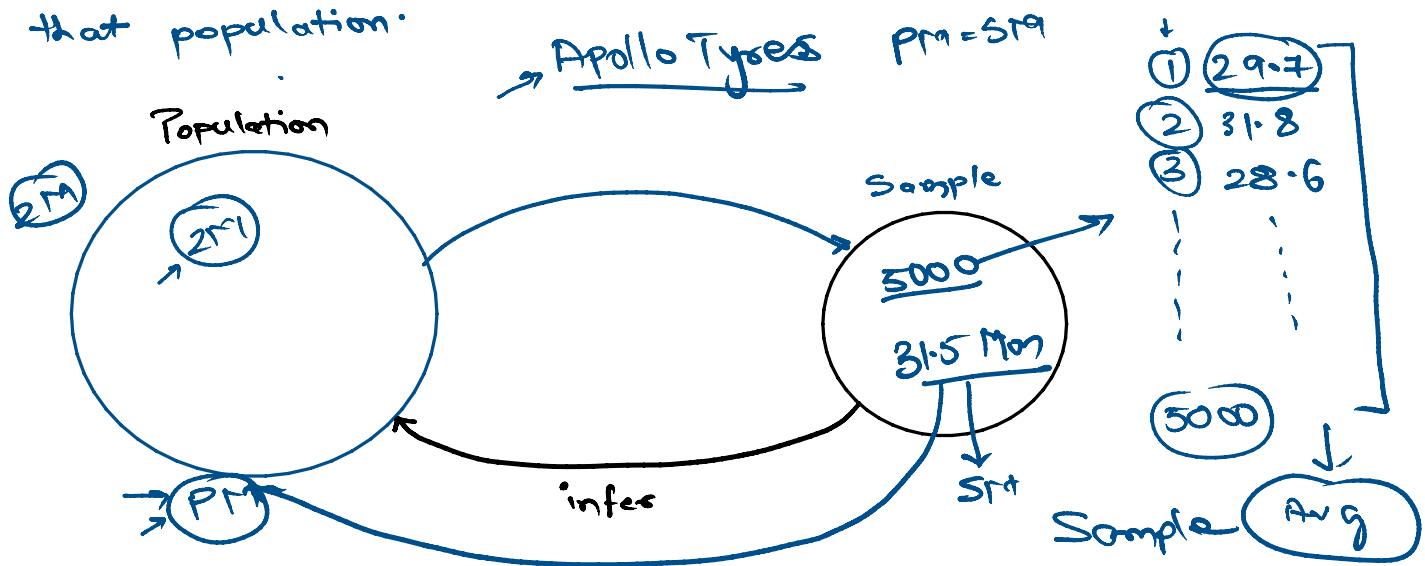
Now we check for this value in Z-score table.

$$z = -1.5 \Rightarrow 0.0668 \text{ or } 6.68\%.$$

Z-score tells us about the percentage of data behind a certain value.

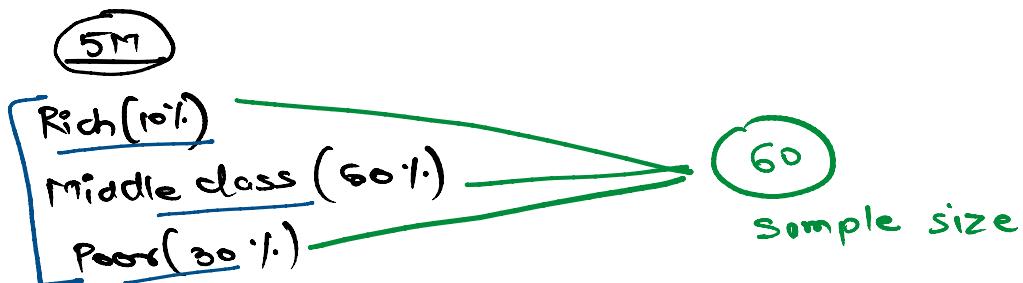
Inferential Statistics →

Involves making predictions or inferences about a population based on the sample data taken from that population.



→ We use sample mean to infer population mean.

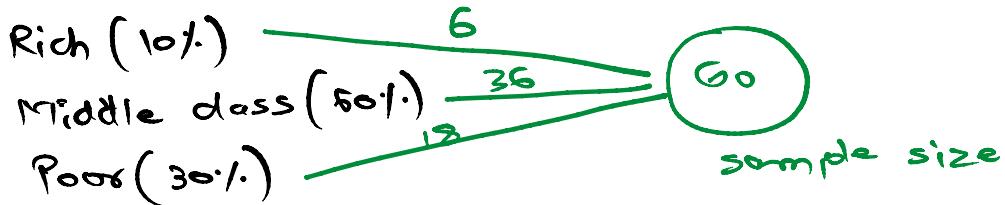
→ The sample should be representative of the population.



Disproportionate Sampling

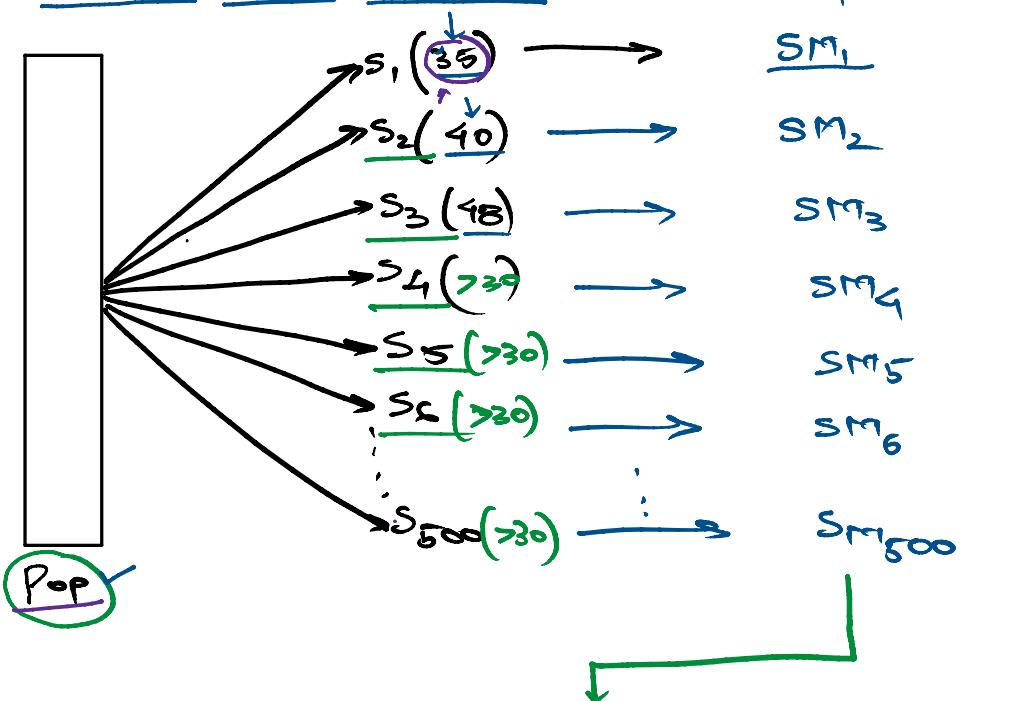
Disproportionate sampling.

5M



Proportionate Sampling.

Central Limit Theorem →



Sample Means

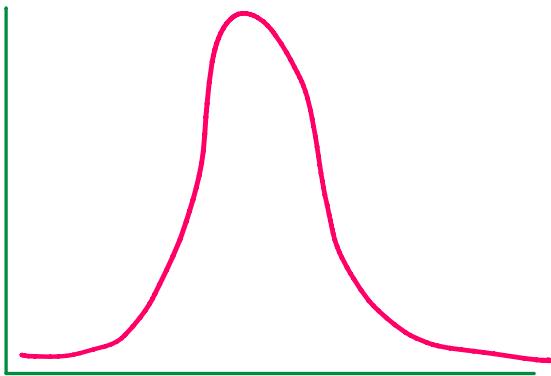
≥ 30
 ≥ 30
 ≥ 200

[]
 $\frac{500 \text{ cm}_s}{}$

Take all the sample means & plot these values.

(sns.kdeplot or sns.distplot)

It will always come out to be a "Normal Distribution".



Note: Sample sizes > 30 .

Takeaway: If we take one sample only from any population, we can safely assume that the sample will be 'Normally Distributed', given that the sample size > 30 .

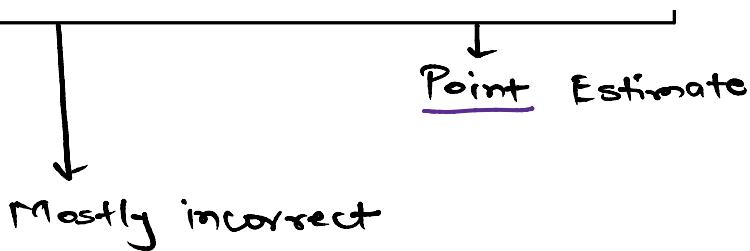
Sample = 50 \rightarrow goes
 ↓
Test these

- ① 24.3 months.
- ② 29.8 months
- ③ 30.5 months.
- ④ 26.9 months.
- ⋮
- ⑤ 31.9 months.

→ Plot these (kdeplot)
 ↓
 Normal Distribution

↓
 Population Mean = Sample Mean

↓
 Point Estimate



With some adjustments, we can say →

$$\text{Population mean} = \text{Sample Mean} \pm \text{Margin}$$

$$= 24.7 \pm 1$$

$$= 24.7 + 1 \quad \text{to} \quad 24.7 - 1$$

$$= 23.7 \quad \text{to} \quad 25.7$$

Interval / Range

90% confidence → $24.7 \pm 1k$

95% confidence → $24.7 \pm 1.5k$

99% confidence → $24.7 \pm \underline{\underline{2k}}$



Confidence Interval.

Formula to calculate the Confidence Interval →

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

Margin

\bar{x} : Sample Mean

s : standard deviation of the sample.

\bar{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% ✓

Ex: Estimate whether the mean lead content in Maggi packets is within the allowed range or not?

≥ 30

Allowed range = 2.5 ppm

$$n = 100$$

$$\bar{x} = 2.3 \text{ ppm}$$

$$S = 0.3 \text{ ppm}$$

$$1 \rightarrow 2.20 \text{ ppm}$$

$$2 \rightarrow 2.43 \text{ ppm}$$

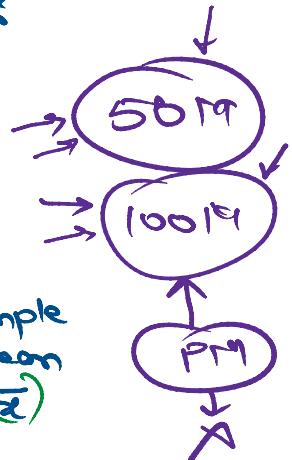
$$3 \rightarrow 2.37 \text{ ppm}$$

$$4 \rightarrow 2.21 \text{ ppm}$$

$$\vdots$$

$$100 \rightarrow 2.28 \text{ ppm}$$

Sample data



$$SM = \frac{2.3 \text{ ppm}}{\sqrt{n}} \times 2.5 \text{ ppm}$$

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

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

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

$$= 2.3 + 0.07 \rightarrow MOE$$

z^*

90%

95%

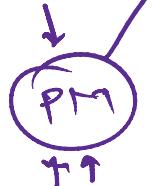
99%

$$= 2.3 \pm 0.07 \rightarrow MOE$$

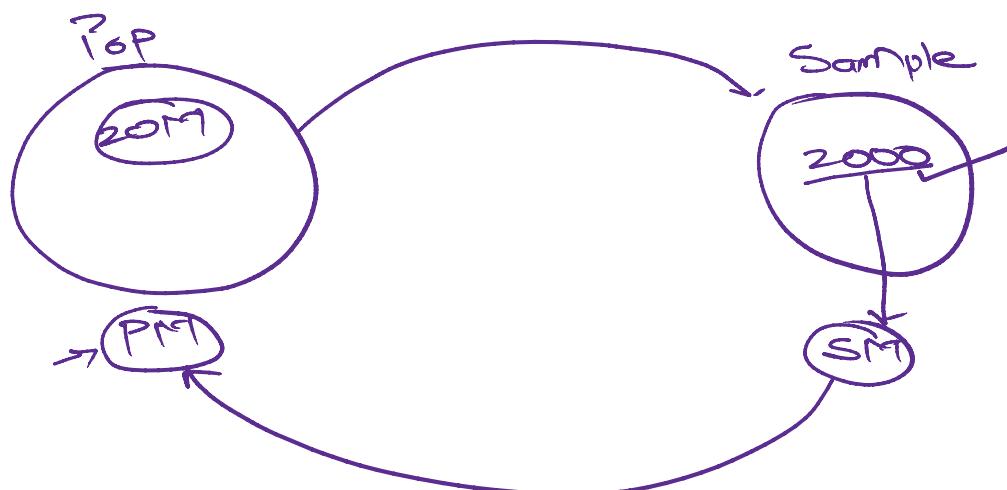
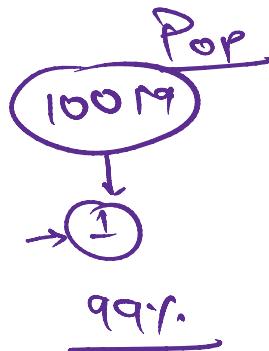
(99%)

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

CI = 2.23 ppm to 2.37 ppm ✓
99% confidence



Maggi



$$PM = SM$$

Mostly incorrect

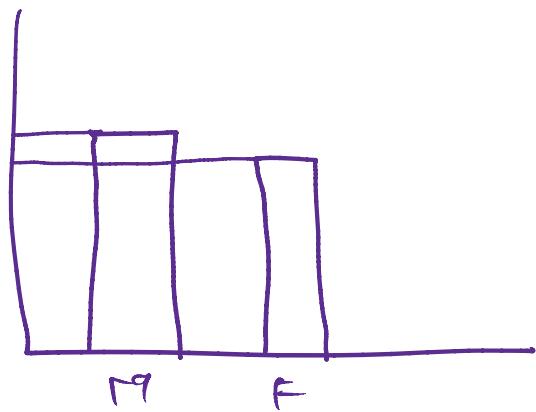
$$PM = SM \pm \text{margin of error}$$

one collection of data

Numerical

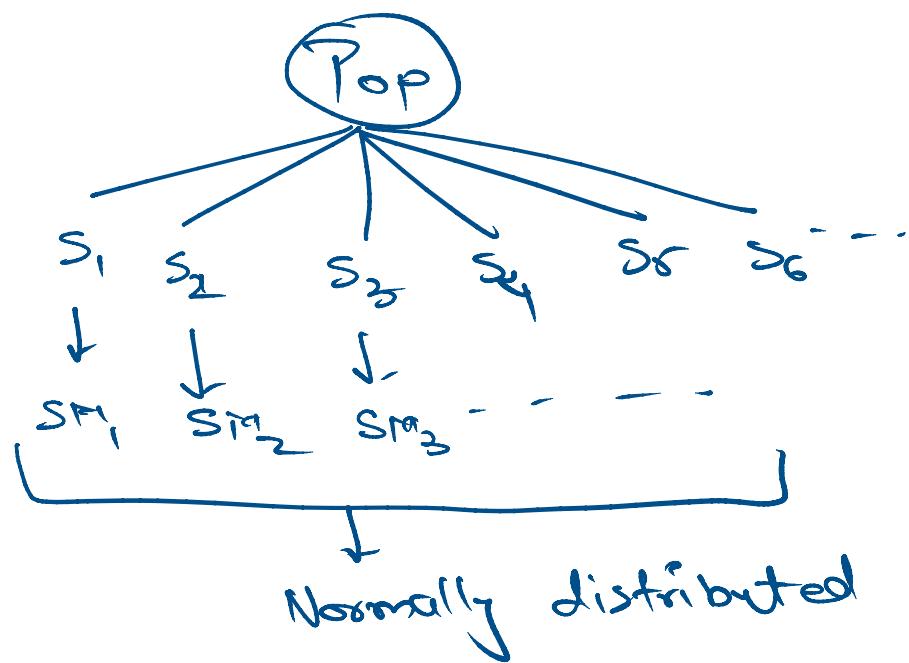
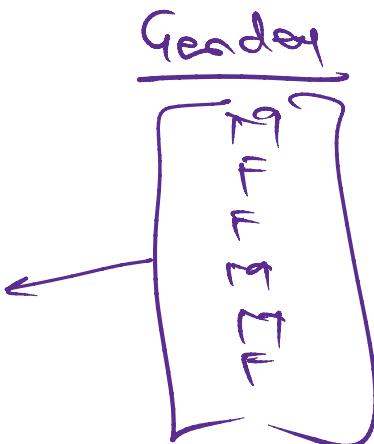
↓

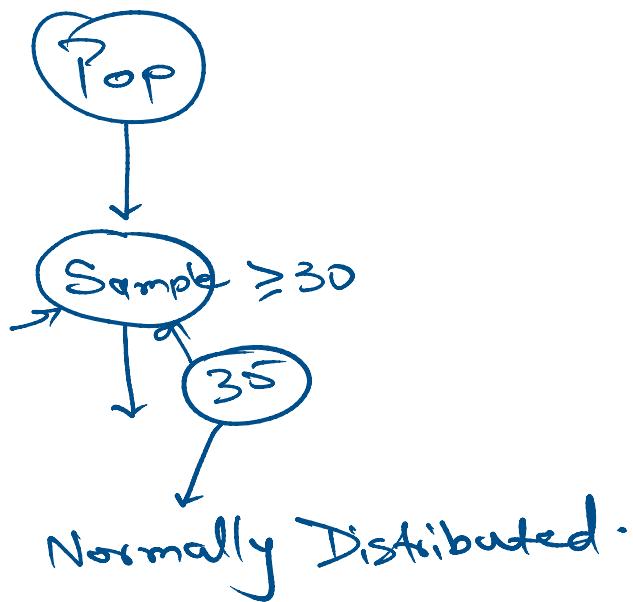
[Histogram]
[Kdeplot]
[Displot]



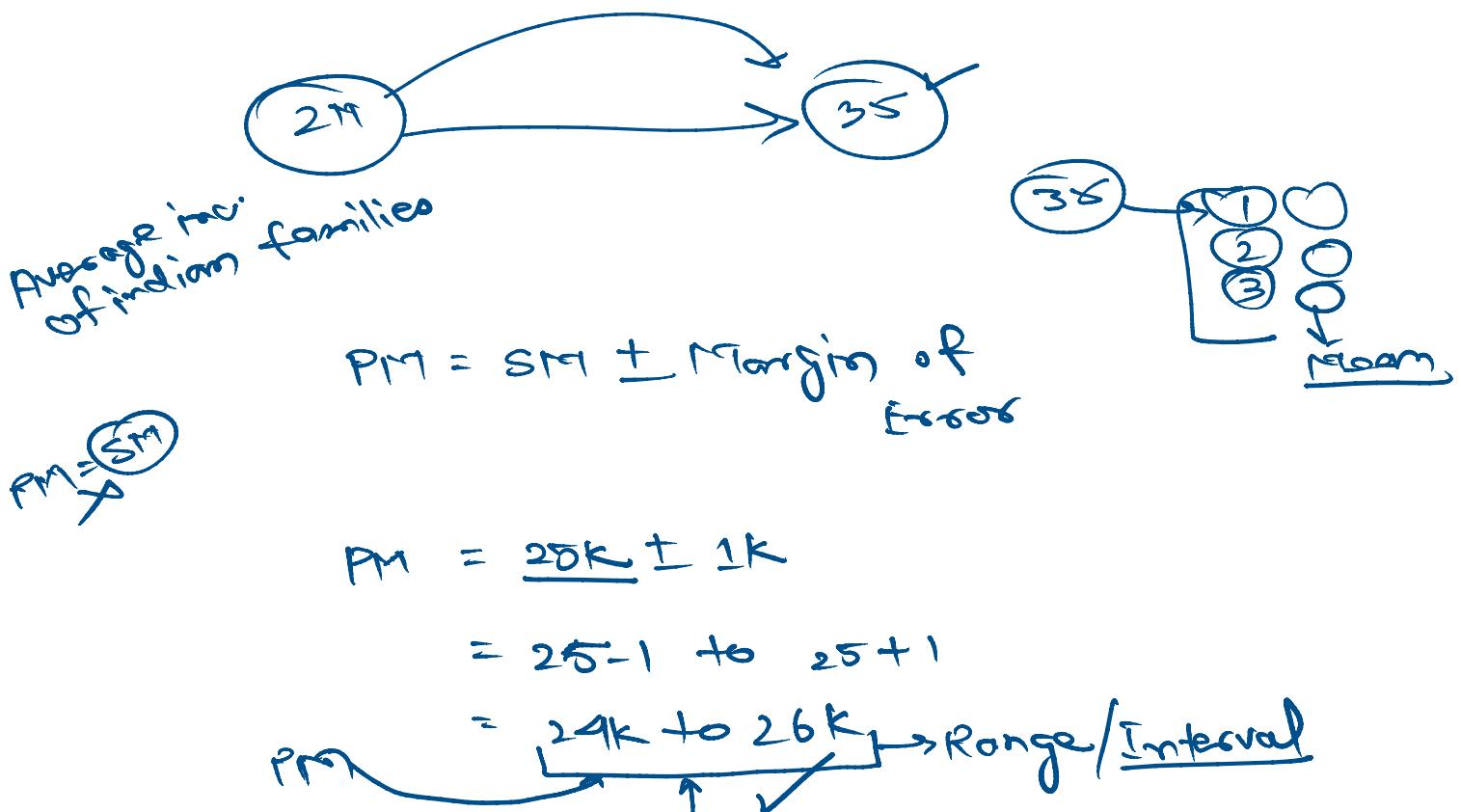
Categories

Bar
Pie



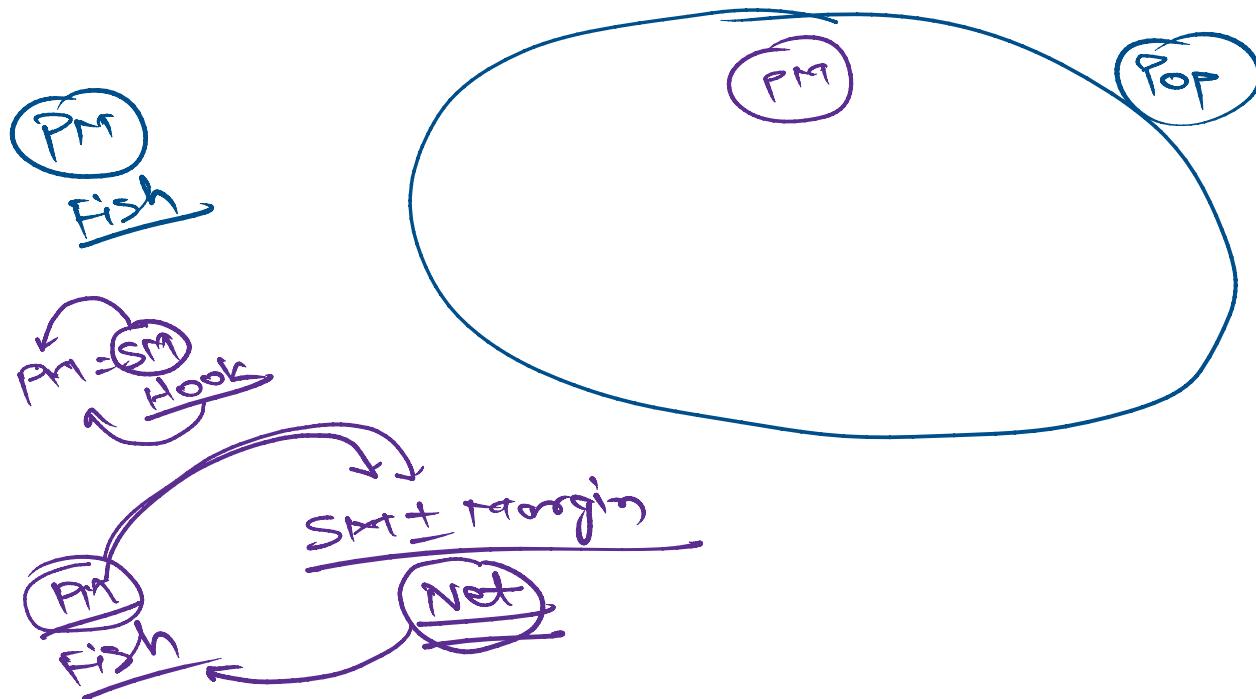


$\geq 30 \rightarrow$ Through CLT
Proven



$$PM = SM \pm \text{Margin}$$

$$= 20K \pm 1K$$



$$\begin{aligned} & SM \pm \text{Margin} \quad \downarrow + \\ \text{90% confidence} & \rightarrow 25K \pm 1K \rightarrow [24K \text{ to } 26K] \\ \text{95% confidence} & \rightarrow 25K \pm 1.5K \rightarrow [23.5K \text{ to } 26.5K] \\ \text{99% confidence} & \rightarrow 25K \pm 2.5K \rightarrow [22.5K \text{ to } 27.5K] \end{aligned}$$

Confidence Interval

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

$$MOI + \frac{\text{Margin}}{\text{Interval}}$$

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

margin / confidence Interval:

$$CI = \bar{x} \pm Z^* \times \frac{s}{\sqrt{n}} \rightarrow \text{Margin of error.}$$

R