

## Inferential Statistics

### Statistical Decision Making Process

Identify the Problem



Collect the data



use descriptive statistics, : visualization , measures.  
Probability

Information



Inferential statistics : Infer about the data

knowledge



Decision



1. Summarization, presentation visualization of the data
2. mean, median, mode  
range, variance, IQR
3. Gain the insights from the data.

1. Involves making predictions or inferences about a larger population based on the observed statistics in the sample data.
2. Involves techniques:  
Hypothesis testing  
Confidence interval

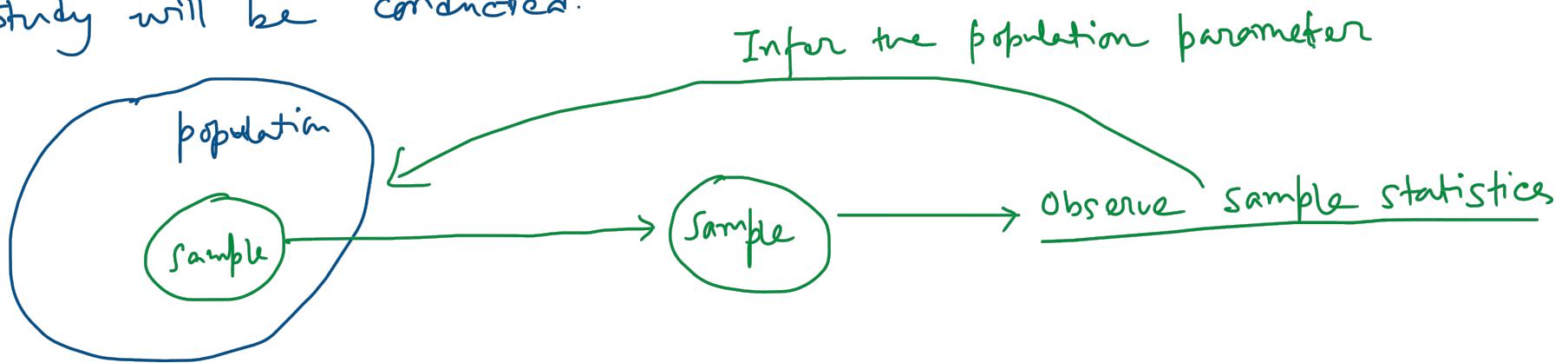
## Population

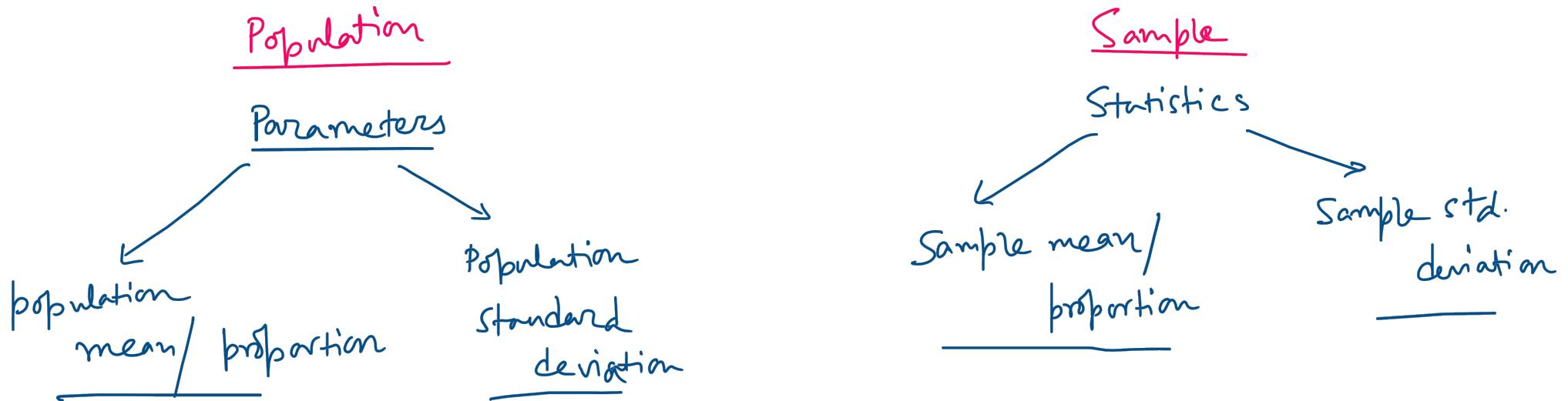
1. The entire group that we want to study

examples: (i) Average height of babies between the age 1 - 5 years in a particular region : [all the babies in that region age between 1 to 5]

## Sample

2. A smaller, representative subset of the population on which the study will be conducted.





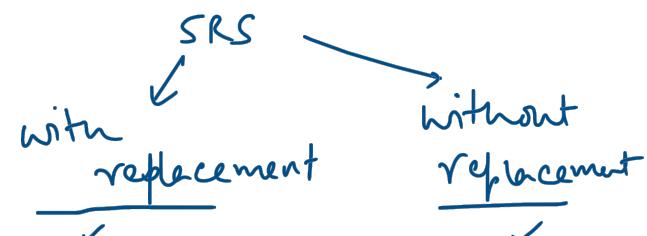
Sample statistics provides an estimate of the population parameters.

There are different sampling techniques

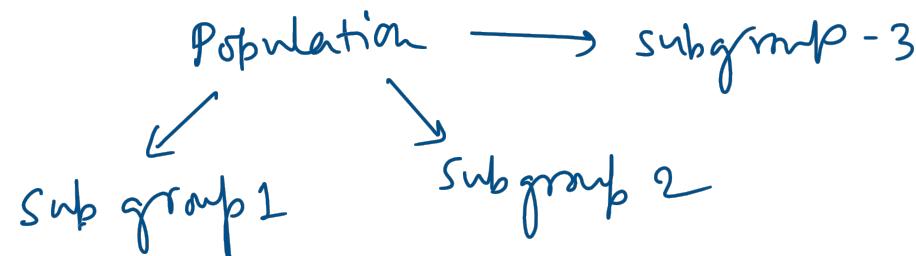
1. Simple Random Sampling (SRS)

$$N = \underline{1000}$$

$\text{rand}(N) \rightarrow 1 \text{ to } 1000$   
 $\underline{100}, \underline{100} \dots$



## 2. Stratified random sampling

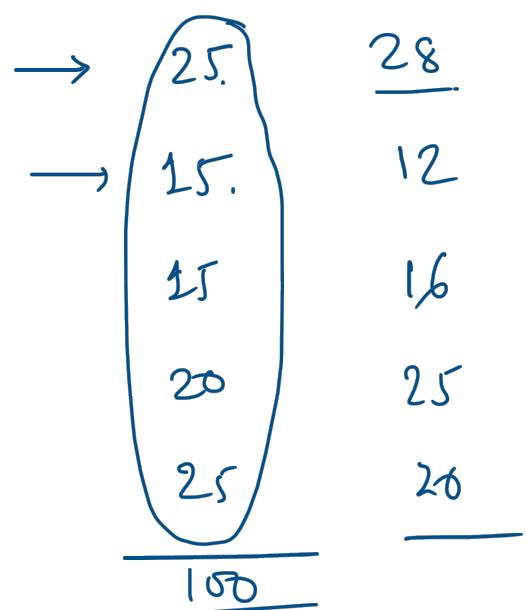


The population is divided into sub groups called strata

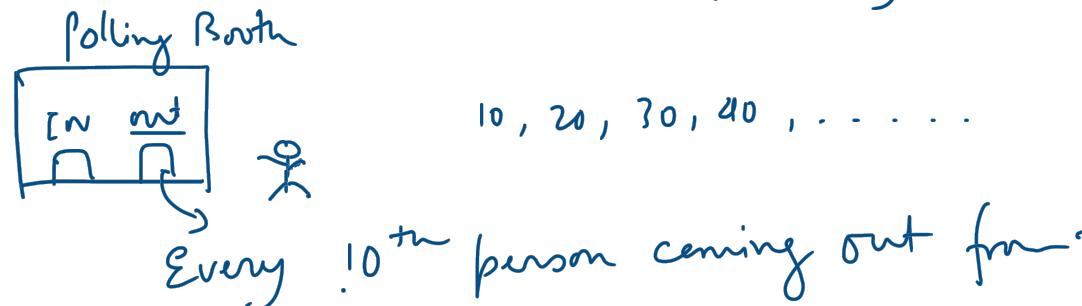
Now from each strata simple random sample are taken in proportion to the population distribution.

Adults in  
. Town

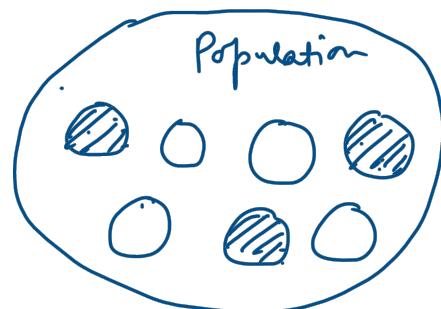
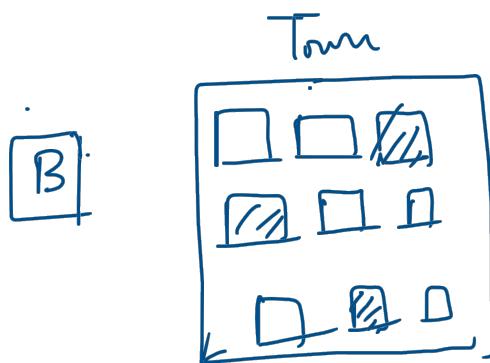
- young adults : 18 yo - 25 yo [25%]
- mid age adults : 25yo - 35yo [15%]
- 35yo - 50 yo [15%]
- 50 yo - 65 yo [20%]
- 65 + yo [25%]



3. Systematic Sampling: Individuals are taken at regular interval from a randomly chosen starting point.



4. Cluster Sampling:



$N$  - clusters.

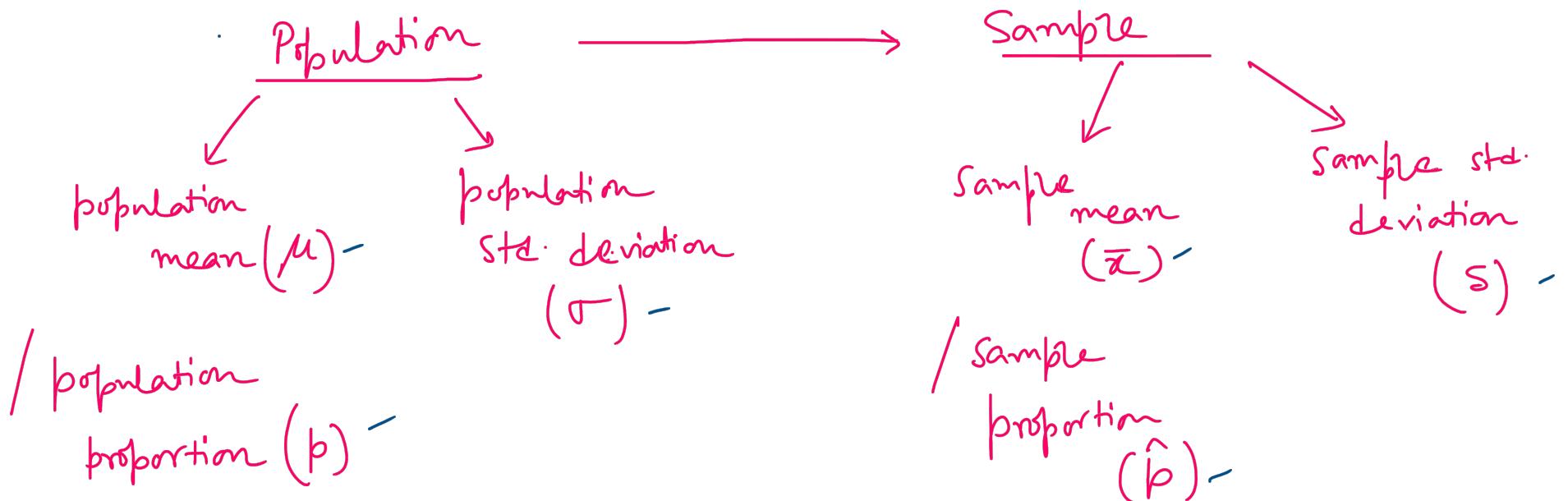
5 cluster out of the  $N$   
cluster

## 5. Convenience Sampling

## 6. Snowball Sampling

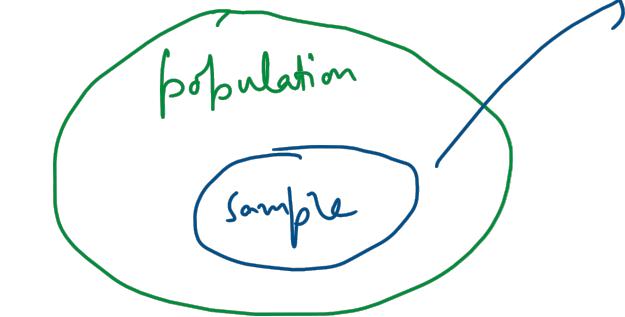


## 7. Purposive Sampling

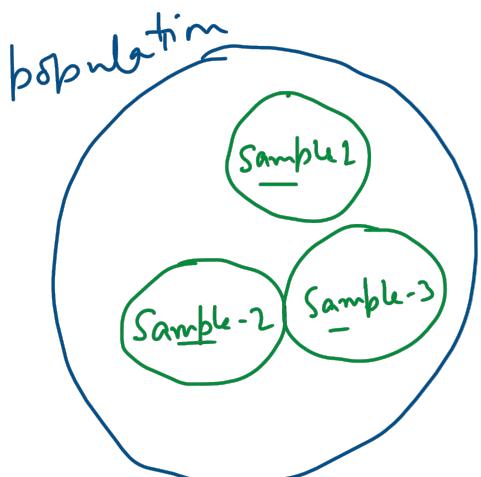


$\bar{x}$  is an estimate of  $\mu$ ,  $\hat{p}$  is an estimate of  $p$   
and  $s$  is an estimate of  $\sigma$ .

## Estimator and Estimates:



Estimator provides an estimate of the certain population parameter solely based on the sample collected..



I want to estimate population mean ( $\mu$ )

I have Sample mean ( $\bar{x}$ ) : an estimator of population mean.

<u>Samples</u>	$\frac{\bar{x}}{\bar{x}^{(1)}}$	:	$s_m \rightarrow \bar{x}^{(n)}$
$s_1$	$\bar{x}^{(1)}$	:	
$s_2$	$\bar{x}^{(2)}$	.	
$s_3$	$\bar{x}^{(3)}$	.	

What properties should an estimator have so that it is called a good estimator?

- No Estimator is perfect
- Suppose an estimator ( $\hat{E}$ ) produces different estimates of different samples taken from the same population

$$\boxed{\hat{E}}$$



$$E[\underline{\epsilon}] = \frac{1}{k} (\epsilon_1 + \epsilon_2 + \dots + \epsilon_k) \approx \underline{\epsilon}$$

$E[\epsilon] = \underline{\epsilon}$  : Then  $\epsilon$  is an unbiased estimator

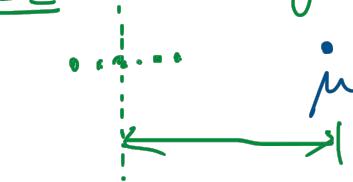
A good estimator on an average gives a very good estimate of the population parameter.

Each estimator is associated with Bias and Variance

$E_1$ : low bias / high variance



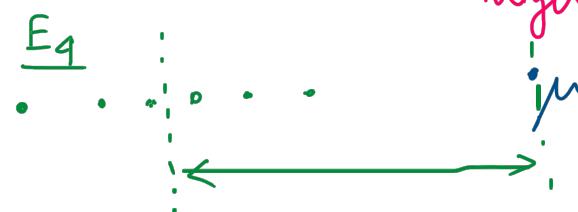
$E_2$ : high bias / low variance



$E_3$



(low bias) / low variance



$$\mathbb{E}[\bar{x}] - \mu : \text{bias}$$

$$\text{Var}(\bar{x}) : \text{variance}$$

## Sample Mean:

Suppose the sample is drawn from the population and sample size is 'n'.

$$x_1, x_2, x_3, \dots, x_n \text{ then } \bar{x} = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$x_i$  are independent and identically distributed random variables with mean = population mean ( $\mu$ ) and std. deviation =  $\sigma$  (pop std. deviation)

$$E[\bar{x}] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n E[x_i]$$

$$x_i \sim (\mu, \sigma^2) \quad i.i.d.$$

$$= \frac{1}{n} \sum_{i=1}^n \mu = \frac{n\mu}{n} = \mu$$

$$E[x_i] = \mu$$

$$\text{Var}(x_i) = \sigma^2$$

$$E[\bar{x}] = \mu$$

Unbiased Estimator of population mean.

$$\text{Var}(\bar{x}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n x_i\right)$$

$$\therefore \text{Var}(\bar{x}) = \frac{1}{n^2} \cdot \sum_{i=1}^n \text{Var}(x_i)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sigma^2$$

$$= \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$$\text{Var}(\bar{x}) = \frac{\sigma^2}{n} \quad \text{and} \quad \text{std}(\bar{x}) = \frac{\sigma}{\sqrt{n}}$$

sample size  $\uparrow \rightarrow \text{variance}(\bar{x}) \downarrow$

$\text{std}(\bar{x}) \rightarrow$  Standard Error

$$x \rightarrow \sigma$$

$$Y = aX$$

$$\text{Var}(Y) = a^2 \text{Var}(X)$$

$$Z = X + Y$$

$$\text{Var}(Z) = \text{Var}(X) + \text{Var}(Y)$$

n: sample size

sample size  $\uparrow \rightarrow \text{standard error} \downarrow$

## Central Limit Theorem :-

$x_1, x_2, \dots, x_n$  n independent and identically distributed random variable with mean  $\mu$  and std. deviation  $\sigma$

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{x} \rightarrow \text{sample mean}$$

as  $n \rightarrow \infty$  the  $\bar{x}_n \rightarrow N(\mu, \frac{\sigma^2}{n})$

As the sample size increases the sample mean approaches towards a normal distribution with mean =  $\mu$  and Std. deviation =  $\frac{\sigma}{\sqrt{n}}$

This is called sampling distribution

The sampling distribution  $\sim N(\mu, \frac{\sigma^2}{n})$  no matter what the original distribution is

## Sample Standard deviation:

$x_1, x_2, \dots, x_n$  : sample values.,  $\bar{x}$  = sample mean

$$= \frac{1}{n} \sum_{i=1}^n x_i$$

Sample Std. deviation  $(S)$  =  $\sqrt{\frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2}$

Why  $(n-1)$ ? and not 'n'?

Population std. deviation  
( $N$  values)

$\mu$ : pop. mean

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

Suppose we define sample std. deviation

$$\hat{s}_n = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}, \text{ then } \hat{s}_n \text{ is a biased estimator of } \sigma$$

always under estimates  $\sigma$ .

Dof : Degree of freedom.

$x_1, x_2, x_3, x_4, x_5$

$$\bar{x} = \frac{x_1 + x_2 + x_3 + x_4 + x_5}{5}$$

if I fix  $\bar{x}$  and change  $x_1, x_2, \dots, x_5$  then how many of these 5 values can we change?

$$x_1 = 2$$

$$\textcircled{1} \bar{x} = \frac{14}{5} = 2.8$$

$$x_2 = 3$$

$$x_3 = 1$$

$$x_4 = 5$$

$$x_5 = 3$$

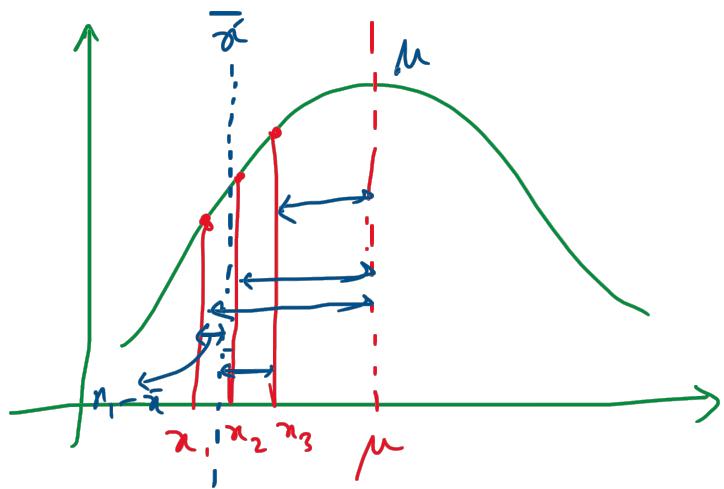
$$\left\{ \begin{array}{l} x_1 = 3 \\ x_2 = 1 \\ x_3 = 5 \\ x_4 = 2 \\ x_5 = 3 \end{array} \right.$$

$$\underline{x_1 + x_2 + x_3 + x_4 + x_5}$$

$$\textcircled{4} \quad \begin{array}{c} n \\ n-1 \\ \vdots \end{array}$$

$$n - \text{dof} = \dots$$

$$\begin{array}{c} \bar{x} \\ \bar{x} \end{array}$$



$$E\left[\sum_{i=1}^n \frac{(x_i - \mu)^2}{n}\right] > \boxed{\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}}$$

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

