

# Sampling Distributions & Sampling Theory

## (i). Population:

The entire set of elements for which the experiments and estimations are to be done.

## (ii) Sample:

A subset of population. All the estimations about the population are made based on the observations of the sample. (also called sample set)

$$\text{Eg } P = \{a, b, c, d, e, f\}$$

$$\text{Then } S_1 = \{a, b, c\}$$

$$S_2 = \{a, d, e\}$$

$$S_3 = \{b, c, f, d\}$$

are examples of samples.

## (iii). The population distribution and parameters:

The population is considered to be a random variable. The distribution of this R.V. is the distribution of the population. The parameters of this distribution are called the parameters of the population.

### Example

- Let  $P = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$   
Let  $S = \{1, 6, 9, 10\}$  be a sample

Now

$$\text{Population mean} = \frac{(1+2+3+4+5+6+7+8+9+10)}{10} \\ = \frac{55}{10} = 5.5$$

$$\text{Sample mean} = \frac{1+6+9+10}{4} \\ = \frac{26}{4} = 6.75.$$

Clearly S.M.  $\neq$  P.M. (This estimation is however close to P.M.)

let  $S_1 = \{1, 10\}$

$$S.M. = 11/2 = 5.5!!$$

But usually we do not know the population mean, thus we don't know which sample mean is correct estimate of P.M.!

The objective of studying sampling theory is to estimate characteristic of the population based on the sample values. Also we need to make formal definitions about the "correctness" of the estimate as we don't know the characteristics of the population.

Notes: (i) Different samples from the same population will result in different estimates ~~also~~ of the parameters. These estimates are themselves themselves treated as R.V. and hence follow some probability distributions, such distributions are called sampling distributions.

(ii). If the distribution of the population is known then only the parameters of this distribution is to be estimated from the observed sample. The confidence in such estimates is to be evaluated.

### (iii) Hypothesis Testing & Goodness of fit

Assume that the distribution of the population is not known. Then we may try to fit some known family of distribution in this fit needs to be evaluated for correctness. This evaluation is called goodness of fit.

However in some cases we may only be interested in testing a hypothesis regarding some property of the population. Based on the observed sample and tests of the sample we may accept or reject the hypothesis.

#### Examples:

- ① Suppose we want to evaluate an algorithm for execution time. Naturally we can't give all possible input (population). So we evaluate on a sample and calculate  $E(T)$  and  $\text{Var}(T)$ ,  $T = \text{execution time}$ . Here the

execution time,  $T$ , is a r.v. Hence  $E(T)$  and  $V(T)$  make sense.

- (2) Jobs come at a call center. The interarrival time is, say, exponentially distributed with parameter  $\lambda$  which is unknown.

Here we know the type of distribution but not the parameter  $\lambda$ . We now estimate  $\lambda$ . For this, we take some samples (i.e. observe the interarrival time for a finite period). Then we obtain an estimate of  $\lambda$  (w/a  $\hat{\lambda}$ ).

Sometimes we are not interested in  $\lambda$  (or so  $\hat{\lambda}$ ) but we are only interested if  $\lambda$  is greater than some threshold value i.e.

$$\lambda = \lambda_0$$

$$\lambda > \lambda_0$$

$$\text{or } \lambda < \lambda_0.$$

If we make an hypothesis that  $\lambda > \lambda_0$ :  $H_1: \lambda > \lambda_0$  ( $\Rightarrow$  more resources needed)

$$H_2: \lambda < \lambda_0$$

$\Rightarrow$  (more resources needed)

If  $H_1$  is accepted then resources sufficient so we proceed to test the hypotheses  $H_1$  &  $H_2$  for acceptance or rejection.

## Parameter Estimation

Let  $P$  be the population. Suppose the probability distribution,  $(X)$ , of  $P$  is known. Only the parameter,  $\theta$ , of  $X$  is not known. Then parameter can be:

$$\theta = \text{Mean of population i.e. } E(X)$$

$$\text{or } \theta = \text{Variance } " " \text{ i.e. } V(X).$$

Procedure:

Collect 'n' experimental outcomes:

$$x_1, x_2, \dots, x_n.$$

Each experimental outcome,  $x_i$ , is a value of a r.v.  $x_i$ . Then the set of random variables

$$x_1, x_2, \dots, x_n$$

is called a sample of size 'n'.

Note that  $x_1, x_2, \dots, x_n$  is a sample and this sample can have so many values from the set  $P$ .

Ex Suppose  $P = \{1, 2, 3, \dots, N\}$ .  
we can have

$$x_1 = 1, \text{ or } 2 \text{ or } 3, \dots \text{ or } N$$

$$x_2$$

$x_1$  can take value  $= 1, 2, \dots, N = x_1$

$x_2$  " " " =  $1, 2, \dots, N = x_2$

$x_n$  " " " =  $1, 2, \dots, N = x_n$

## Random Sample

The set of random variables  $x_1, x_2, \dots, x_n$  is called a random sample of size 'n' from the population,  $P$ , with distribution of the R.V.  $X$  as  $F(x)$ , provided that  $x_1, x_2, \dots, x_n$  are pairwise independent with some distribution  $f^n$  i.e.

$$f_i(x) = f(x) \quad \forall i \in \mathbb{N}$$

This holds if the samples are drawn without replacement. If replacement is not done the outcome of R.V.  $x_2$  will be affected by the outcome of  $x_1$  & so on. In that case

$$P(x_1 = x_1, x_2 = x_2, \dots, x_n = x_n)$$

$$= \frac{1}{N} \cdot \frac{1}{(N-1)} \cdots \frac{1}{N-n+1}$$

$$= \frac{(N-n)!}{N!}$$

## Statistic

Any function  $T(x_1, x_2, \dots, x_n)$  of the observations  $x_1, x_2, \dots, x_n$  is called statistic since it is a f<sup>n</sup> of r.v., a statistic is itself a R.V. and its distribution is

called The sampling distribution of  $T$ .

## (i) Sample Mean:

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

Although ' $\bar{x}$ ' is the distribution of the population, we write  $\bar{x}$  as sample mean to understand that it represents joint PDF  $f(x_1, x_2, \dots, x_n)$ .

## (ii) Sample Variance:

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

## Estimator:

Let  $\theta$  be a parameter of the population. If its statistic (i.e.  $\bar{x}$  of  $x_1, x_2, \dots, x_n$ ) is the estimate of  $\theta$ , then  $\hat{\theta}$  is called estimator of  $\theta$ . i.e.  $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$  is called the estimator of  $\theta$ . A specific value of the estimator  $\hat{\theta}$  is called an estimate of  $\theta$ .

## Unbiased Estimator

A statistic  $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$  is called an unbiased estimator of  $\theta$  if

$$E[\hat{\theta}(x_1, x_2, \dots, x_n)] = \theta$$

## Estimate

Let  $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$  be an estimator of  $\theta$ . A specific value of  $\hat{\theta}$  obtained by substituting observations of  $x_1, x_2, \dots, x_n$  is called estimate of  $\theta$ .

$$\underline{\theta} = \underline{\theta}(x_1, x_2, \dots, x_n)$$

Ex

The sample mean is an unbiased estimator of the population mean  $\mu$ .

For  $E[\bar{X}] = \frac{1}{n} E \left[ \sum_{i=1}^n X_i \right]$

$$= \frac{1}{n} \sum_{i=1}^n E[X_i]$$

$$= \frac{1}{n} [E[X_1] + E[X_2] + \dots + E[X_n]]$$

We know that DF of each  $X_i$  is ~~same~~ the DF of  $X$  i.e.

$$F_{X_i}(x) = F(x) \text{ for } i = 1, 2, \dots$$
  
$$\therefore E[X_i] = E[X]$$

$$\therefore E[\bar{X}] = \frac{1}{n} [\mu + \mu + \dots + \mu]$$

$$= \frac{1}{n} [n\mu]$$

$$= \mu$$

~~∴~~ This estimator is an unbiased estimator of population mean.

Computing Variance of the estimator sample mean:

$$V[\bar{X}] = \sum_{i=1}^n V_{\text{am}} \left[ \frac{X_i}{n} \right]$$

Note Consider Rolling of die 3 times

$X_1$  = first outcome,  $X_2 = 2^{\text{nd}}$  &  $X_3 = 3^{\text{rd}}$  outcome

$$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n \left[ \frac{3}{2} X_i \right] = \frac{1}{n} \sum_{i=1}^n E[X_i]$$

$$= \frac{1}{n} [E[X_1] + E[X_2] + E[X_3]]$$

$$= \frac{1}{n} [3 \cdot 5 + 3 \cdot 5 + 3 \cdot 5] = 3 \cdot 5 = \mu$$

$$= \text{Var} \left[ \sum_{i=1}^n \frac{x_i}{n} \right]$$

We know for independent RV  $x_1$  &  $x_2$   
 ~~$\text{Var}[x_1 + x_2] = \text{Var}(x_1) + \text{Var}(x_2)$~~  ?

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

~~$$= \text{Var} \left[ \frac{1}{n} \sum_{i=1}^n x_i \right]$$~~

~~$$= \frac{1}{n} V \left[ \sum_{i=1}^n x_i \right]$$~~

~~$$= \frac{1}{n} [V(x_1) + V(x_2) + \dots + V(x_n)]$$~~

~~$$= \frac{n}{n^2} \left[ \sum_{i=1}^n V[x_i] \right]$$~~

For independent RV  $x_1$  &  $x_2$ , we have

$$V[\alpha_1 x_1 + \alpha_2 x_2] = \alpha_1^2 V[x_1] + \alpha_2^2 V[x_2]$$

Generalizing and applying here

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

$$= \frac{1}{n^2} \sum_{i=1}^n V[x_i]$$

$$= \frac{1}{n^2} \sum_{i=1}^n V[x_i]$$

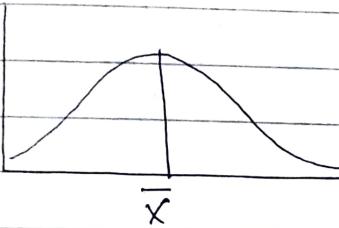
$$= \frac{1}{n^2} [n\sigma^2]$$

$$\Rightarrow E[V[\bar{X}]] = \sigma^2/n$$

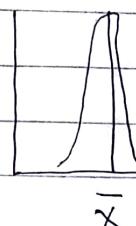
This means that  $\bar{X}$  as estimate of  $\mu$  will have less variance when  $n$  is large. Typically

$$V[\bar{X}] = 0$$

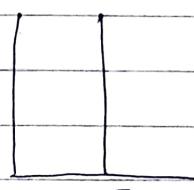
$n \rightarrow \infty$



large  $V[\bar{X}]$



small  $V[\bar{X}]$



$V[\bar{X}] \rightarrow 0$

$\therefore$  As  $n \rightarrow \infty$ , The accuracy of  $\bar{X}$  as an estimate of population mean is excellent.

E.S.

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

is an unbiased estimator of the population variance.

Note: Be careful that  $E[\bar{X}] = \mu$  and not  $\bar{X} = \mu$ . This means taking many samples & finding  $\bar{X}$  will lead to different  $\bar{X}$  (However); the average of all  $\bar{X}$  will be  $\mu$ .

Proof:

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

$$= \frac{1}{(n-1)} \sum_{i=1}^n [x_i^2 - 2x_i\bar{x} + \bar{x}^2]$$

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

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

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

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

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

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

$$\text{But } E[x_i^2] = V[x_i] + (E[x_i])^2$$

$$= \sigma^2 + \mu^2 \quad \text{--- } ③$$

$$\& \mathbb{E}[\bar{x}^2] = \mathbb{V}[\bar{x}] + (\mathbb{E}[\bar{x}])^2$$

$$= \frac{\sigma^2}{n} + \mu^2 \quad (4)$$

Substituting (3) & (4) in (2)

$$\mathbb{E}[s^2] = \frac{1}{(n-1)} \sum_{i=1}^n (\sigma^2 + \mu^2) - \frac{n}{(n-1)} \left[ \frac{\sigma^2}{n} + \mu^2 \right]$$

$$= \frac{n}{(n-1)} (\sigma^2 + \mu^2) - \frac{n}{(n-1)} \left( \frac{\sigma^2}{n} + \mu^2 \right)$$

$$= \frac{n}{(n-1)} \left[ \sigma^2 + \mu^2 - \frac{\sigma^2}{n} - \mu^2 \right]$$

$$= \frac{n}{(n-1)} \left[ \frac{n\sigma^2 - \sigma^2}{n} \right]$$

$$= \frac{n}{(n-1)} \left[ \frac{(n-1)\sigma^2}{n} \right]$$

$$\Rightarrow \boxed{\mathbb{E}[s^2] = \sigma^2}$$

$\therefore s^2$  is an unbiased estimator of the population variance  $\sigma^2$ .

Ex

Let the population be "heights of 20 students in a class" normally distributed with

$$\mu = 175 \text{ cm}$$

$$\text{ & } \sigma = 5 \text{ cm. } \Rightarrow \sigma^2 = 25$$

$$P = \{163, 159.5, 159.14, \dots, 186.41, 178.99, 185\}$$

Let us calculate a sample mean  $\bar{x}$  with sample size = 3.

$$\text{Let } S_1 = \{x_1, x_2, x_3\}$$

$$S_2 = \{x_1, x_2, x_3\}$$

There are other samples possible.

$$\text{Sample mean for } S_1 = \bar{x}_1 = \frac{163 + 169.5 + 159.14}{3}$$

$$\Rightarrow \boxed{\bar{x}_1 = 160.5 \text{ cm.}}$$

$$\text{Sample mean for } S_2 = \bar{x}_2 = \frac{186.41 + 178.99 + 185.2}{3}$$

$$\Rightarrow \boxed{\bar{x}_2 = 183.5 \text{ cm}}$$

Sampling distribution of  $\bar{x}$

$$\underbrace{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n}_{X}$$

Number of samples  $\neq$  Sample size.

We know for normal distribution

$$E[\bar{X}] = \mu \text{ & } E(\bar{X}) = \frac{\sigma^2}{n}$$

$$\text{Var}[\bar{X}] = \frac{\sigma^2}{n}$$

If the population is  $N(\mu, \sigma^2)$

then distribution of  $\bar{X}$  is  $N(\mu, \frac{\sigma^2}{n})$

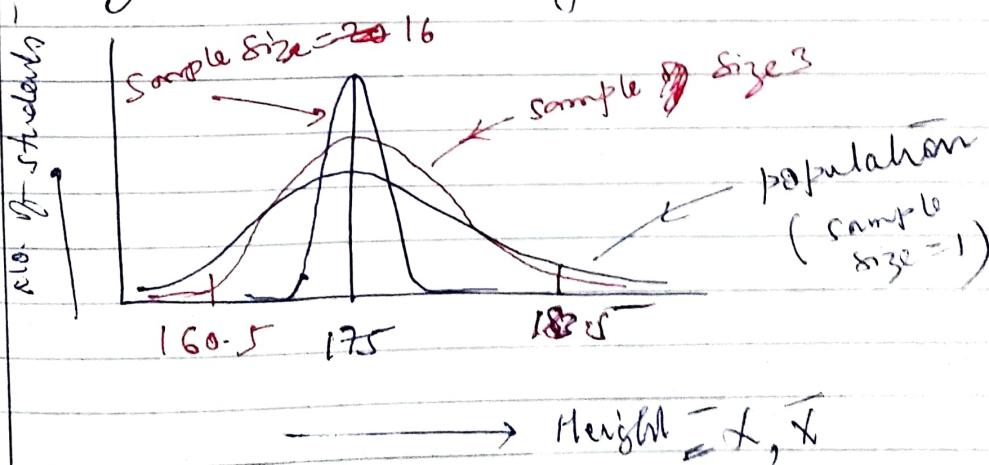
If  $X$  represents height of class then

$$\mu_{\bar{X}} = 175 \text{ cm} \text{ & } \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} = 2.9 \text{ cm}$$

$\therefore \bar{X}$  has distribution ~~for~~  $N(\mu, \frac{\sigma^2}{n})$

Now we see Sample  $\frac{\text{mean}}{\text{size}}$  Variance =  $\frac{\sigma^2}{n} = \frac{\sigma^2}{n}$

As sample size increases  $\leftarrow n \rightarrow 20$   
 then ~~sample~~:  $V[\bar{X}] \rightarrow 0$  i.e.  $\mu_{\bar{X}} = \mu = 175$   
 which is obvious because there is  
 only one sample of size 20.



Note Population distribution is a special case of sampling dist with  $n=1$