

# Sampling

1. Our aim is to draw a genuine conclusion about a large group of objects/individuals (population). Instead of examining the entire group (may be difficult/impossible) we only inspect only a small part of this population, called a sample. The process of obtaining sample is sampling.

i) Suppose we would like to <sup>draw</sup> conclusions about the heights of 10k students (pop<sup>n</sup>) by examining only, say 100 students (sample) selected from the population.

ii) We may wish to draw conclusions about the colours of 1k marbles (pop<sup>n</sup>) in a bag, (sample) selecting a sample of 20 marbles from the bag (each marble selected is returned after its colour is observed  $\rightarrow$  NR)

2. i) "Population"  $\rightarrow$  used to denote the observations/measurements (rather than individuals/objects)

ii) Population  $\rightarrow$  finite/infinite. The no. is called the size of population ( $N$ ). The no. in the sample ( $n$ ) is called the sample size.

iii)  $N = 10k, 1k$ .  $n = 100, 20$ . (above)

iv) Suppose we would like to draw conclusions about the fairness (unbiased) of a specific coin by tossing it repeatedly. Population  $\rightarrow$  all possible tosses of the coin ( $N$ : infinite)

A sample  $\rightarrow$  may be obtained by examining, say 70 tosses of the coin and noting the % of heads/tails ( $n: 70$ )

### 3. Sampling $\rightarrow$ WR, WOR

i) After drawing ~~an~~ an object from a bag we may or may not replace same into the bag before we draw again.

a particular object comes up repeatedly  $\rightarrow$  the object comes up once

(WR)

(WOR)

ii) A finite popl<sup>n</sup> that is sampled with replacement can theoretically be considered infinite since samples of any size can be drawn without exhausting the population.

### 4. How to choose a sample?

One way to do this for finite populations is to make sure that each member of the popl<sup>n</sup> has equal chance of being included in the sample. — which is then called a random sample.

As inference from sample to popl<sup>n</sup> can't be certain, we use the language of probability in any statement of conclusions.

### 5. Parameter (of population)

Population <sup>known</sup>  $\rightarrow$  if its prob dist<sup>n</sup>  $f(x)$  of the associated r.v.  $X$  is known.

If  $X \sim$  binomial dist<sup>n</sup> we say that the popl<sup>n</sup> is binomially distributed or that we have a binomial popl<sup>n</sup>.

There will be certain quantities that appear in  $f(x)$ , e.g.  $\mu$ ,  $\sigma$  (for normal)  $\mu$ ,  $\sigma$  (for normal), other quantities

moments, skewness etc. can be determined in terms of  $\mu$ ,  $\sigma$  etc. All such quantities are of ten called population parameters.

For a given / known population  $f(x)$ , population parameters are also known.

If prob. distr.  $f(x)$  of the population is not known precisely, although we may have some idea of, or at least be able to make some hypothesis concerning, the general behaviour of  $f(x)$ .

Ex: We may have some reason to suppose that a particular population is normally distributed. In that case we may not know one or both of the values  $\mu$  and  $\sigma$  and so we might wish to draw statistical inferences about them.

## 6. Statistics (of sample)

$X$ : values are various heights

$X_1$ : height of the 1st individual,  $x_1$  be its value

$X_2$ : height of the 2nd individual,  $x_2$  be its value.

In general, a sample of size  $n$  described by the values  $x_1, x_2, \dots, x_n$  of the random variables  $X_1, X_2, \dots, X_n$

SRSWR:  $X_1, X_2, \dots, X_n \rightarrow$  iid having prob. distr.  $f(x)$ . Their joint pdf  $P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = f(x_1) f(x_2) \dots f(x_n)$

## Sample Statistic or Statistic

Any quantity obtained from a sample for the purpose of estimating a popl<sup>n</sup> parameter.

A statistic for a sample of size  $n$  can be defined as a function of  $n$  r.v.  $X_1, X_2, \dots, X_n$  i.e.,  $g(X_1, X_2, \dots, X_n) \rightarrow$  another random variable

Its values are represented by  $g(x_1, x_2, \dots, x_n)$ .  
Corresponding to each popl<sup>n</sup> parameter,  $\theta$  a statistic to be computed from the sample.

One of the imp. prob<sup>s</sup> of sampling is to decide how to form the statistic that will best estimate a given popl<sup>n</sup> parameter.

popl<sup>n</sup> parameters:  $\mu, \sigma$

sample statistics:  $m, s$  etc.

## 7. Sampling Distribution

The prob<sup>b</sup> dist<sup>n</sup> of a sample statistic  $g(X_1, X_2, \dots, X_n)$  is called the sampling dist<sup>n</sup> of the statistic.

Alt. Consider all possible samples of size  $n$  to be drawn from the popl<sup>n</sup>. For each sample we compute the statistic.

In this way, we obtain the dist<sup>n</sup> of the statistic called sampling dist<sup>n</sup>.

Mainly, we are interested to compute sampling dist<sup>n</sup> of sample mean, sample variance etc.

## 8. Sampling Dist<sup>n</sup> of Sample Mean

$X_1, X_2, \dots, X_n \rightarrow$  iid r.v. for a random sample of size  $n$ .

$$\bar{X} = \frac{1}{n} \sum X_j \rightarrow \text{sample mean.}$$

If  $x_1, x_2, \dots, x_n$  denote values obtained in a particular sample of size  $n$ , then the mean for that sample  $\bar{x} = \frac{1}{n} \sum x_j$ .

Let  $f(x)$ : prob. dist<sup>n</sup> of some given  $\mu$  and variance  $\sigma^2$  popl<sup>n</sup> with mean  $\mu$  and variance  $\sigma^2$ , from which we draw a sample of size  $n$ .

i) the mean of the sampling dist<sup>n</sup> of means, denoted by  $\mu_{\bar{x}} = E(\bar{x}) = \mu$ , the popl<sup>n</sup> mean

ii) If the popl<sup>n</sup> is infinite and the sampling is random or if the popl<sup>n</sup> is finite and sampling is with replacement,

the variance of the sampling dist<sup>n</sup> of means, denoted by  $\sigma_{\bar{x}}^2 = E[\bar{x} - \mu]^2 = \frac{\sigma^2}{n}$ ,  $\sigma$ : the popl<sup>n</sup> variance,

iii) If the popl<sup>n</sup> is finite ( $N$ ) and the sampling is without replacement (WOR), then

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \left( \frac{N-n}{N-1} \right) \rightarrow \frac{\sigma^2}{n}, \text{ as } N \rightarrow \infty$$

iv) If the popl<sup>n</sup>  $X \sim N(\mu, \sigma^2)$ , then the sample mean  $\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ .

v) If the popl<sup>n</sup> from which samples are taken, which has a prob. dist<sup>n</sup> with mean  $\mu$  and variance  $\sigma^2$  (not necessarily normal distribution), then

$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$  is asymptotically normal,

i.e.,  $\lim_{n \rightarrow \infty} P(Z \leq z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$

(Popl<sup>n</sup> is infinite / sampling is WR)

## 9. Sampling Dist<sup>n</sup> of Proportions

Pop<sup>l<sup>n</sup></sup>  $\rightarrow$  infinite & binomially distrib. <sup>prob. that any given member exhibits</sup> but  
Consider all possible samples of size 'n' drawn from this pop<sup>l<sup>n</sup></sup> and for each sample, we determine the statistic that is  
the proportion P of successes

P: proportion of heads turning up in n tosses

Sampling dist<sup>n</sup> of Proportions (Whose mean  $\mu_p$  and s.d.  $\sigma_p$ ) as

$$\mu_p = p, \quad \sigma_p = \sqrt{\frac{pq}{n}} = \sqrt{\frac{p(1-p)}{n}}$$

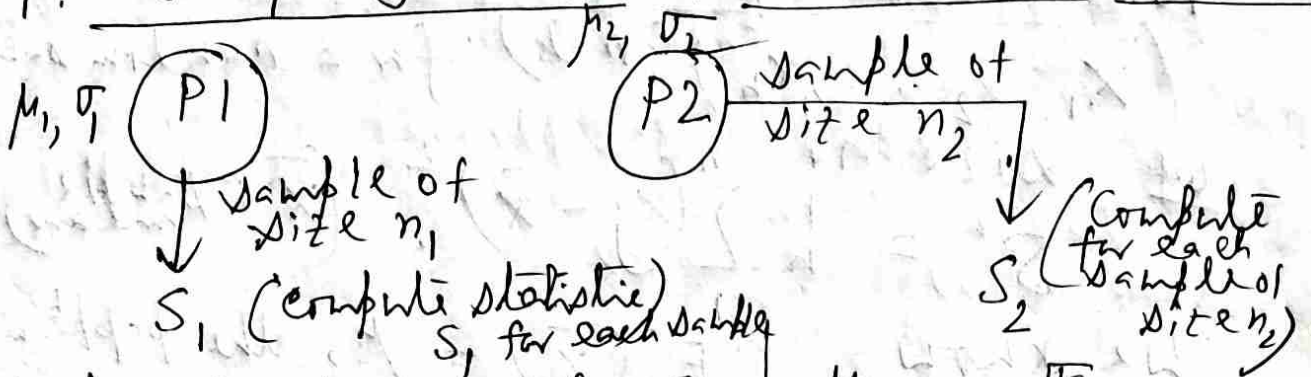
$$E(\bar{X}) = \mu, \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

[  $\mu_p$  &  $\sigma_p$  are obtained by putting  $\mu = p$  and  $\sigma = \sqrt{pq}$  ]

For large n ( $> 30$ ), the sampling dist<sup>n</sup> is close to normal dist<sup>n</sup>.  
For finite pop<sup>l<sup>n</sup></sup> (sampling is

WOR),  $\sigma_p \leftarrow \frac{\sigma}{\sqrt{n}}$

# 9. Sampling Dist<sup>n</sup> of Differences & Sums



Let mean & s.d. of  $S_1$  be  $M_{S_1}$  and  $\sigma_{S_1}$   $M_{S_2}, \sigma_{S_2}$

Dist<sup>n</sup> of  $(S_1 - S_2)$   $\rightarrow$  sampling dist<sup>n</sup> of differences of the statistics

Mean:  $M_{S_1 - S_2} = M_{S_1} - M_{S_2}$  ( $S_1, S_2$  are indep)

S.D.:  $\sigma_{S_1 - S_2} = \sqrt{\sigma_{S_1}^2 + \sigma_{S_2}^2}$

If in particular  $S_1$  &  $S_2$  be sample means of two populations, denoted by  $\bar{X}_1, \bar{X}_2$  then

$$M_{\bar{X}_1 - \bar{X}_2} = M_{\bar{X}_1} - M_{\bar{X}_2} = \mu_1 - \mu_2$$

$$\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Results hold also for finite pop<sup>n</sup> if sampling is WR. (sampling dist<sup>n</sup> of difference of means)

The standardized variable

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \approx \text{Normally distributed if } n_1, n_2 \text{ be large } (n_1, n_2 \geq 30)$$