

Sampling of size  $n$  from population

- model: slips drawn from bag

- drawing

with replacement (i.i.d.)

without replacement (SRS)

} very little difference if  $n \ll$  size of population

if  $\frac{n}{\text{size of population}} \leq 0.1$

Think of independence as events not influencing one another.

Extend idea of independence to variables  $X, Y$ :

say they are independent if they are not associated

Several facts about quantitative vars.

1. If  $X, Y$  are quant. vars and  $X$  has a mean  $\mu_x$ ,  $Y$  has a mean  $\mu_y$ , then their sum  $X+Y$  has mean  $\mu_x + \mu_y$ .

Their difference  $X-Y$  has mean  $\mu_x - \mu_y$ .

Ex-) If you're a golfer and average 91 strokes on course A and 83 strokes on course B. If you play both courses on one day, and take

$X =$  your score on course A

$Y =$  " " " " B

then

average sum  $X+Y$  will be  $\mu_x + \mu_y = 91 + 83 = 174$

" difference  $X-Y$  " "  $\mu_x - \mu_y = 91 - 83 = 8$ .

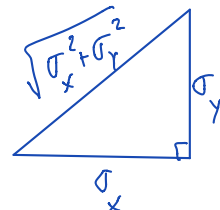
2. If  $X, Y$  are independent and  $X$  has a s.d.  $\sigma_x$ ,  $Y$  has a standard deviation  $\sigma_y$ , then both

$X+Y$

$X-Y$

} have standard deviation

$$\sqrt{\sigma_x^2 + \sigma_y^2}$$



Ex.]

Say Laura is a bowler whose scores have

$$\text{mean } \mu = 161$$

$$\text{s.d. } \sigma = 14$$

Laura bowls two games adding her scores

$$X_1 = \text{score in 1st}$$

$$X_2 = \text{" " 2nd}$$

Know

$$S = X_1 + X_2 \text{ has mean } 161 + 161 = 322$$

$$\text{" s.d. } \sqrt{14^2 + 14^2} = \sqrt{2(14)^2} = 14\sqrt{2} \\ = \sigma\sqrt{2}.$$

If she bowls 3 games, then her summed / total score

$$X_1 + X_2 + X_3 \text{ has mean } = 3\mu = 3(161)$$

$$\text{s.d.} = \sqrt{14^2 + 14^2 + 14^2} = 14\sqrt{3}.$$

3. Corollary to 2:

If  $X_1, X_2, \dots, X_n$  is an i.i.d. sample from a population with mean  $\mu$ , s.d.  $\sigma$ , then the

a) mean for the sum:  $n\mu$ , s.d. for sum:  $\sigma\sqrt{n}$

b) mean for  $\frac{1}{n}(X_1 + X_2 + \dots + X_n)$  (i.e., the mean for the sample mean  $\bar{X}$ )

$$\mu_{\bar{X}} = \mu \quad \left( \text{sample means are unbiased estimators of the population mean} \right)$$

and the s.d. of  $\bar{X}$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}.$$

Ex.]

If Laura bowls 3 games, expect her average of the three

$$\bar{X} = \frac{1}{3}(X_1 + X_2 + X_3)$$

to have a distribution with mean 161 and s.d.  $\frac{14}{\sqrt{3}}$

---

If she bowls 20 games, then her average

$$\bar{X} = \frac{1}{20}(X_1 + X_2 + \dots + X_{20})$$

will have mean = 161 and s.d.  $\frac{14}{\sqrt{20}}$ .

Central Limit Theorem:

1) The sum of  $X_1, \dots, X_n$  (i.i.d. sample from quantitative population) has an approximate normal distribution for large enough  $n$  (and given what we learned above, that normal dist. will have mean  $n\mu_x$  and s.d.  $\sigma_x \sqrt{n}$ ).

2) The average  $\frac{1}{n}(X_1 + \dots + X_n)$  will likewise be approximately normal for large enough  $n$  (with mean  $\mu_x$ , s.d. =  $\frac{\sigma_x}{\sqrt{n}}$ ).

Both claims assume i.i.d. samples  $X_1, \dots, X_n$  taken from a population. They also apply for SRS's, if rule of thumb (sampling less than 10% of population).

Stat 145, Mon 18-Oct-2021 -- Mon 18-Oct-2021  
Biostatistics  
Spring 2021

-----  
Monday, October 18th 2021  
-----

Wk 8, Mo  
Topic:: Central Limit Theorem  
Read:: Lock5 5.2

Variables can

- have an association, or
- not have an association.

We also talk about independent variables, which is roughly  
the same as

Examples:

1. If we draw twice from a bag and take

$X$  = 1st outcome

$Y$  = 2nd outcome

then  $X$  and  $Y$  are

- i) independent if sampling "with replacement"  
call this an i.i.d. random sample of size 2

- ii) approximately independent if the composition of the bag is  
little changed after the first draw

2. If we draw  $n$  times from a bag and take

$X_1$  = 1st outcome

$X_2$  = 2nd outcome

.

.

.

$X_n$  = nth outcome

the  $X_i$  are

- i) independent if sampling "with replacement"  
call this an i.i.d. random sample of size  $n$

- ii) approximately independent if the composition of the bag is little changed after by the draws

Rule of thumb: size  $n$  of sample is  $\leq 10\%$  of size of bag's contents

A random variable  $X$  is one that is numeric for each case

- sex: F/M we think of as categorical (binary)
- $X(\text{case}) = 0$  if case=female, 1 if case=male is a random variable (turns outcomes into numbers)

Some facts about independent normal random variables

- If  $X$  and  $Y$  are independent normal random variables, with
  - $X \sim \text{Norm}(\mu_1, \sigma_1)$  — means "has a normal dist. with mean  $\mu_1$ , sd  $\sigma_1$ "
  - $Y \sim \text{Norm}(\mu_2, \sigma_2)$
- then  $X+Y$  (their sum) is  $\sim \text{Norm}(\mu_1 + \mu_2, \sqrt{\sigma_1^2 + \sigma_2^2})$
- then  $X-Y$  (their difference) is  $\sim \text{Norm}(\mu_1 - \mu_2, \sqrt{\sigma_1^2 + \sigma_2^2})$

Ex.: Suppose Ray and Joan are bowlers. Their scores have normal dists

$$R \sim \text{Norm}(142, 17)$$

$$J \sim \text{Norm}(138, 22)$$

$$R+J \sim \text{Norm}(280, 27.803)$$

How likely is it for them, in one game, to have a combined score  $> 350$ ?

Answer comes from  $1 - \text{pnorm}(350, 280, 27.803)$

- If we draw an i.i.d. random sample of size  $n$ , each  $X_i \sim \text{Norm}(\mu, \sigma)$ , then the
  - sum =  $X_1 + \dots + X_n$  is  $\text{Norm}(n \mu, \sigma \sqrt{n})$
  - avg =  $(X_1 + \dots + X_n) / n$  is  $\text{Norm}(\mu, \sigma / \sqrt{n})$

Central Limit Theorem

Suppose a random sample of size  $n$  is drawn from the population either

- with replacement (so it is i.i.d.), or
- with  $n$  smaller than 10% of the full population.

If the variable of interest is quantitative and  $n$  is large enough, then

the sum  $X_1 + \dots + X_n$  is approximately normal

the mean  $(X_1 + \dots + X_n)/n$  is approximately normal

If the variable of interest is binary categorical and  $n$  is large enough, then the sample proportion has approximately a normal distribution.

Explorations using apps at

[https://onlinestatbook.com/stat\\_sim/sampling\\_dist/index.html](https://onlinestatbook.com/stat_sim/sampling_dist/index.html)

<https://shiny.calvin.edu:3838/scofield/samplingDists/>

<https://shiny.calvin.edu:3838/scofield/cltProportions/>

## Central Limit Theorem

In summary, here is the take-away from the **Central Limit Theorem**.

---

Suppose you have a random sample of size  $n$  that is either

- i.i.d., or
- an SRS, with the sample size  $n$  being no more than 10% of the size of the population.

In the case that

1. the variable under consideration is quantitative, having population mean  $\mu$  and standard deviation  $\sigma$ , then the sampling distribution for the sample mean  $\bar{X}$  is approximately  $\text{Norm}(\mu, \sigma / \sqrt{n})$  for  $n$  large enough.
  2. the variable under consideration is binary categorical, having population proportion  $p$ , then the sampling distribution for the sample proportion  $\hat{p}$  is approximately  $\text{Norm}(p, \sqrt{p(1-p)/n})$  for  $n$  large enough.
-

Since

- null distributions
- randomization distributions
- bootstrap distributions

are all specialized versions of sampling distributions, then so long as the sample statistic in question is the sample's *mean*  $\bar{X}$  or the sample *proportion*  $\hat{p}$ , we can expect the CLT to apply to these as well.