

## 6 - Sampling Distribution of the Mean

We are about to establish one of the central links between *descriptive statistics* and *probability theory*. This will lead us to *inferential statistics*. The key idea is this:

**The sample mean  $\bar{X}$  is a random variable.**

**Example** Suppose we have a class of  $N = 100$  students. Each student has a certain age,  $X$ . Assume that all the values of  $X$  are:

20	20	20	20	23	20	18	20	21	24
20	23	20	22	23	19	22	21	19	22
19	19	21	24	19	21	22	25	21	21
23	21	20	18	20	19	20	21	20	23
20	19	21	21	22	21	22	19	22	21
23	20	21	19	22	24	22	21	21	22
21	23	20	19	20	25	20	22	22	24
20	19	20	22	24	21	19	20	20	20
22	20	19	21	19	23	21	25	21	22
19	22	20	20	23	23	20	20	20	26

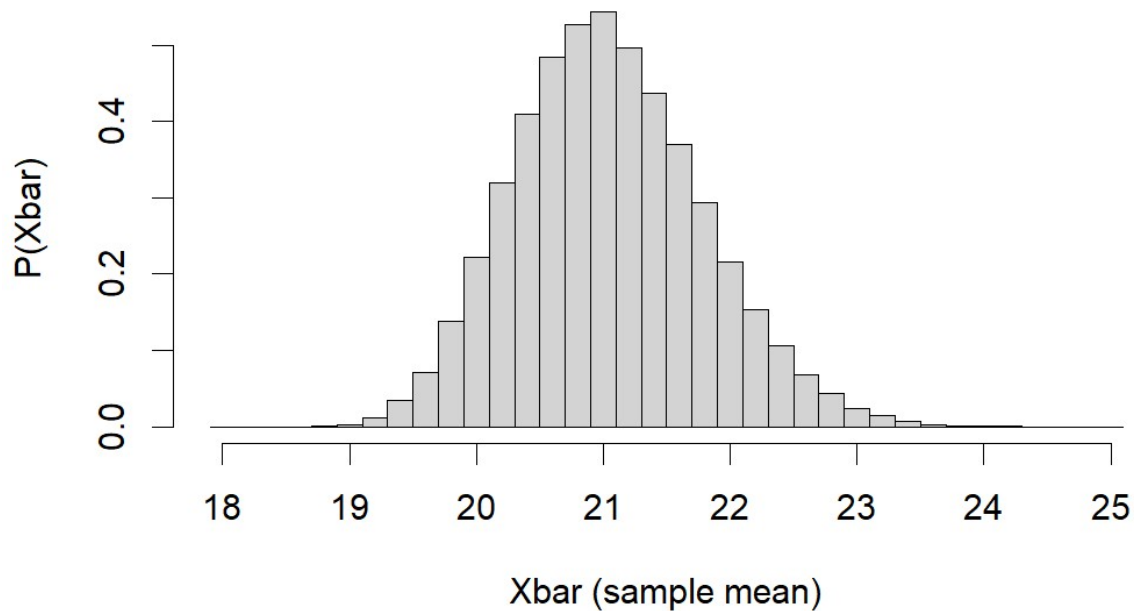
Using R, we find that the population mean and standard deviation are:



Suppose we randomly select a sample of  $n = 5$  students.

It turns out that the sample mean  $\bar{X}$  is a random variable whose distribution is/has:

## Sampling Distribution of $\bar{X}$ ( $n = 5$ )



**Definition** The distribution of  $\bar{X}$  is also called the *sampling distribution of the mean*.

### Observations:

- The values of  $\bar{X}$  are now decimal values (like 20.8) instead of integer values like  $X$ .
- Therefore, there are more *possible* values of  $\bar{X}$ . (i.e., 19.8, 20.0, 20.2, ...)
- If we increased the sample size  $n$ , then  $\bar{X}$  would become a *continuous* variable.
- It does not make much difference if we use *replacement* or *no replacement* as long as the population size  $N$  is much larger than the sample size  $n$ .

### Summary of Symbols

$X$  = the random variable (both sample and population)

$N$  = population size

$\mu$  = population mean value of  $X$

$\sigma$  = population std dev. of  $X$

$n$  = sample size

$\bar{X}$  = sample mean of  $X$

$s$  = sample std dev of  $X$

*Sampling Distribution of the Mean – In General*

Suppose  $X$  is any random variable for individuals selected from a population of size  $N$ . Then  $X$  has a probability distribution with parameters:

- mean  $\mu$
- std. dev.  $\sigma$

If samples of size  $n$  are randomly selected from the population, then the sample mean  $\bar{X}$  is a random variable with its own probability distribution called the *sampling distribution of the mean*. The parameters of  $\bar{X}$  are:

- mean  $\mu_{\bar{X}} = \mu$
- std. dev.  $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$

Think of  $\mu_{\bar{X}}$  as “the average of the averages.”

Think of  $\sigma_{\bar{X}}$  as the typical error between  $\bar{X}$  and  $\mu$ . (the “standard error of the mean.”)

## Central Limit Theorem

Using more advanced concepts from probability theory, it is possible to prove:

**Central Limit Theorem** Suppose a random variable  $X$  is defined on a population of size  $N$ . If a sample of  $n$  individuals are selected independently from the population, then the sampling distribution of  $\bar{X}$  is:

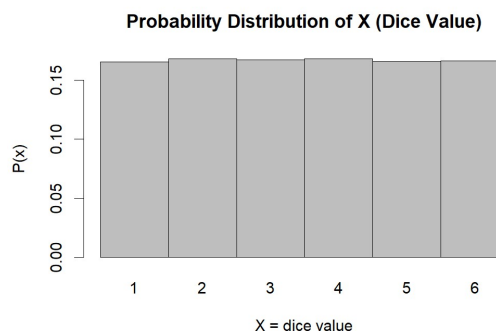
- *normal* if  $X$  was normal
- *approximately normal* even if  $X$  was not normal
  - the approximation becomes better and better as  $n \rightarrow \infty$

*As a practical rule, if  $n \geq 30$  then  $\bar{X}$  follows a normal distribution (even if  $X$  does not).*

**Example (Dice Rolls)** Suppose you roll  $n = 1$  six-sided die. Let  $X$  = the number that turns up. Then the probability distribution of  $X$  is given by:



$x$	$P(X = x)$
1	$1/6$
2	$1/6$
3	$1/6$
4	$1/6$
5	$1/6$
6	$1/6$



The mean and standard deviation of  $X$  are:

$$\mu = \sum [x \cdot P(x)] = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6} = 3.5$$

$$\begin{aligned} \sigma^2 = \sum [(x - \mu)^2 \cdot P(x)] &= (1 - 3.5)^2 \cdot \frac{1}{6} + (2 - 3.5)^2 \cdot \frac{1}{6} + (3 - 3.5)^2 \cdot \frac{1}{6} + \\ &\quad (4 - 3.5)^2 \cdot \frac{1}{6} + (5 - 3.5)^2 \cdot \frac{1}{6} + (6 - 3.5)^2 \cdot \frac{1}{6} = 2.91666 \dots \end{aligned}$$

$$\sigma =$$

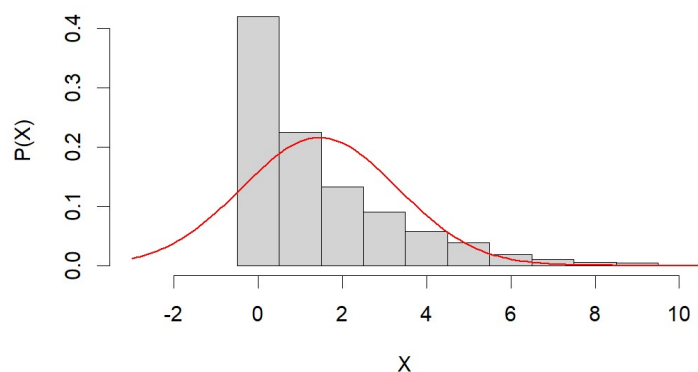
What is the sampling distribution of  $\bar{X}$  if we roll  $n = 5$  dice?

What is the sampling distribution of  $\bar{X}$  if we roll  $n = 20$  dice?

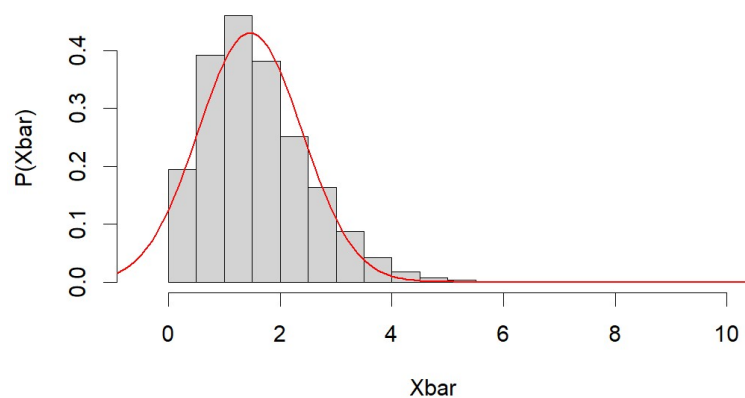
What is the sampling distribution of  $\bar{X}$  if we roll  $n = 100$  dice?

**Example** If the distribution of  $X$  is highly *skewed* then it requires a larger  $n$  before the sampling distribution of  $\bar{X}$  becomes close to normal.

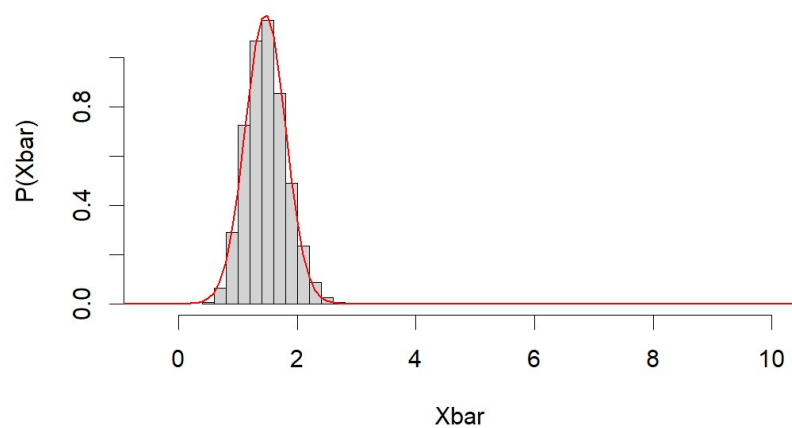
**Distribution of  $X$  (highly skewed)**



**Sampling Distribution of  $\bar{X}$  ( $n=4$ )**



**Sampling Distribution of  $\bar{X}$  ( $n=30$ )**



**Example** In human engineering and product design, it is important to consider the weights of people. Assume that the population of male BCIT students has normally distributed weights, with mean 173.2 lbs and a standard deviation of 29.5 lbs.



- a. Find the probability that if a male student is randomly selected, his weight is greater than 200 lbs.
- b. An elevator has a maximum weight capacity of 7200 lbs. Find the probability that 36 randomly selected male students will exceed the elevator's weight capacity.

**Example** Assume that the population of adult body temperatures has a mean  $37.0^\circ\text{C}$ . Also assume that the population standard deviation is  $0.62^\circ\text{C}$ .

- a. If a sample of  $n = 108$  is randomly selected, find the probability of getting a mean of  $36.8^\circ\text{C}$  or lower.

- b. Suppose you take a sample of  $n = 108$  randomly selected adults and find that the mean of their body temperatures is  $36.8^\circ\text{C}$ . Would this be *unusually low* (in a statistical sense), given the data about the population of adult body temperatures? What does your result suggest?



### *Finite Population Correction Factor*

We have used the formula

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

under the assumption that individuals were sampled *independently* from the population, which requires either:

- the population is infinite or practically infinite ( $n < 5\%$  of  $N$ ), or
- sampling is *with replacement*

If sampling is done *without replacement* where  $n \geq 5\%$  of  $N$ , then we need to adjust the formula above:

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

The expression  $\sqrt{\frac{N-n}{N-1}}$  is called a *finite correction factor*.

**Example** Suppose we collect midterm scores for all 150 students in MATH 3042. We find the mean is 71% and the standard deviation is 6.5%. If 40 students are chosen without replacement, find the probability that their mean score is greater than 73.

