

# Sampling Distribution for a Mean

# Topics

1.  $\mu$  vs.  $\bar{x}$
2. Sampling distribution for a mean
3. Central limit theorem

我们来看一下sample mean的分布到底是什么样子的

# $\mu$ vs. $\bar{x}$

$\mu$  = population mean

- Parameter

$\bar{x}$  = sample mean

- Statistic
- Estimator of  $\mu$

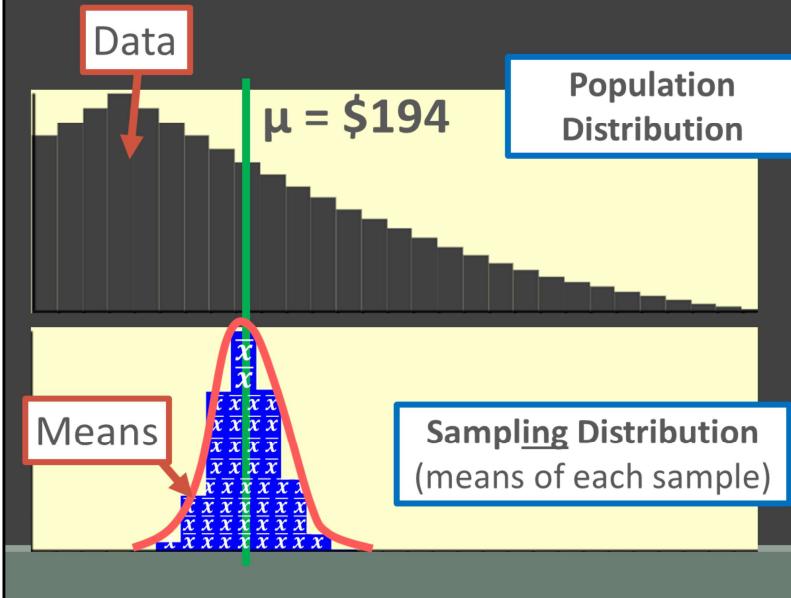
# Topics

1.  $\mu$  vs.  $\bar{x}$

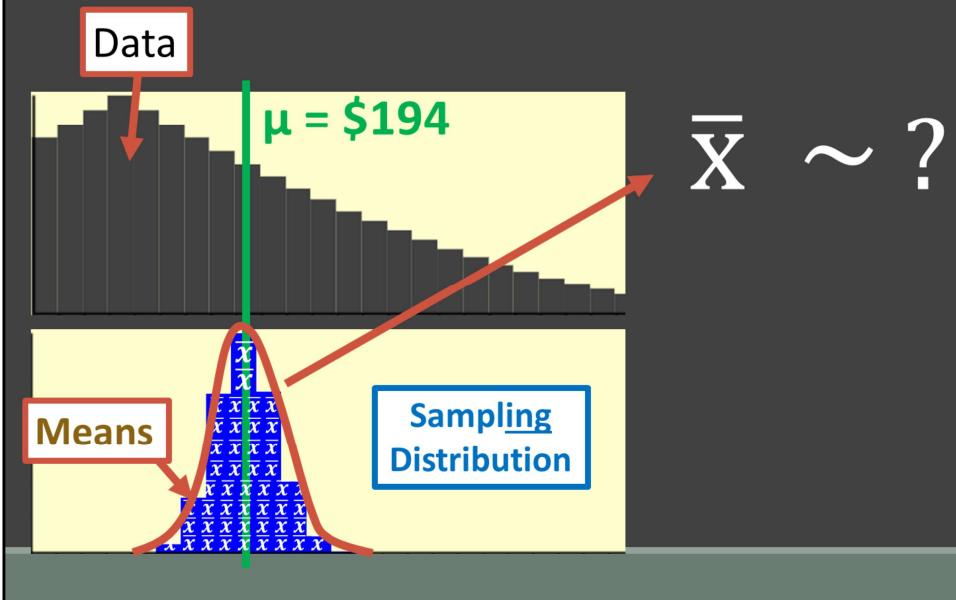
**2. Sampling distribution for a mean**

3. Central limit theorem

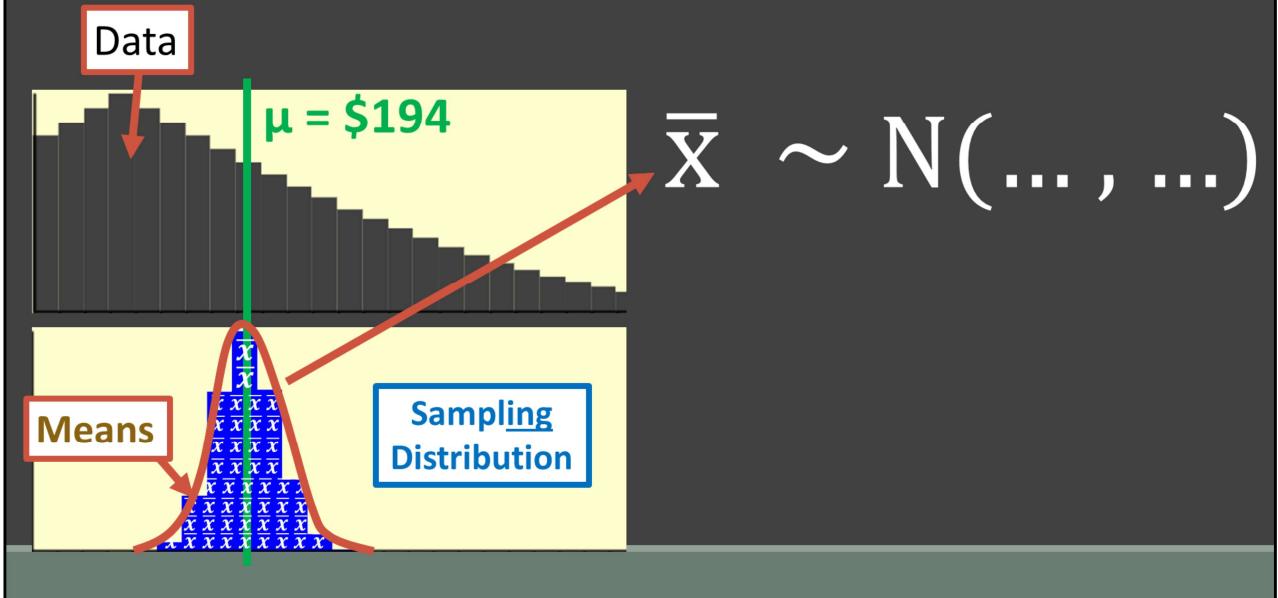
# Sampling Dist. for a Mean



## Sampling Dist. for a Mean



## Sampling Dist. for a Mean



我们发现都比较像Normal，事实也是如此，满足一定条件的情况下，sample mean的分布就趋近于Normal Distribution

那这里也是需要我们来确认一下两个系数： $\bar{x}$ 的期望和方差  
先来看一下它的定义

$$\text{Population Mean: } \mu = \frac{x_1 + x_2 + \dots + x_N}{N}$$

$$\text{Estimator (Statistic): } \bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

这里面，对于 $\bar{X}$ 来说我们回想一下得到均值的过程

我们是不是从population里面随机抽取一个个体，记录它的值，比如就是这里的 $x_1$ ， $x_1$ 是不是从population中随机抽取得到的？

那是不是 $x_1$ 的取值有很多种可能性？

那它是什么？RV！

那同理， $x_2$ 到 $x_n$ 其实都是随机变量，如果我们每次抽取都是with replacement，那每个随机变量还都是独立的

$$\text{Population Mean: } \mu = \frac{x_1 + x_2 + \dots + x_N}{N}$$

$$\text{Estimator (Statistic): } \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

Independent Random Variables:

$X_1, X_2, \dots, X_n \sim \text{Distribution of } X \text{ with } (\mu, \sigma^2)$

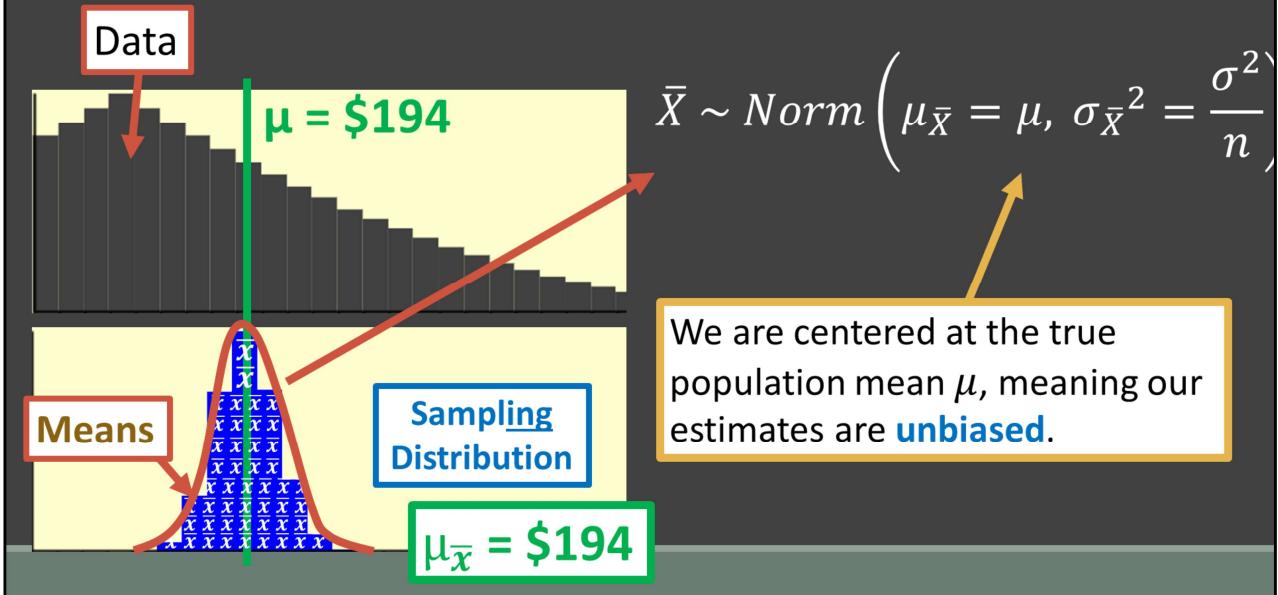
所以我们可以把xbar改写成:  $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$

其中,  $X_i$ 是独立且拥有共同分布的随机变量, 他们都服从于population Distribution, 也就是他们的期望都是mu

那xbar的均值是? Mu

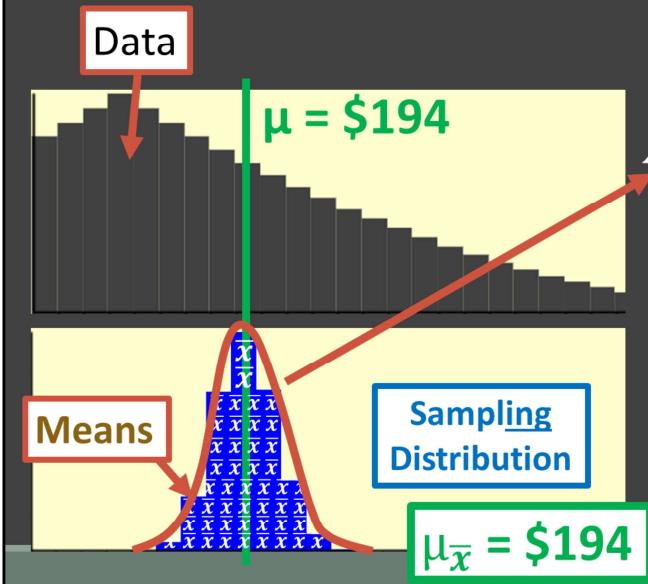
方差是? Sigma^2/n

## Sampling Dist. for a Mean



所以我们得到了 $\bar{x}$ 的分布

## Sampling Dist. for a M



$\sigma \rightarrow$  The higher the spread in the population data, the higher the spread in the sampling distribution.

$$\bar{X} \sim \text{Norm} \left( \mu_{\bar{X}} = \mu, \sigma_{\bar{X}}^2 = \frac{\sigma^2}{n} \right)$$

We **divide by n**, so spread affected by sample size.  
Higher sample size means **less variation**.

# Topics

1.  $\mu$  vs.  $\bar{x}$
2. Sampling distribution for a mean
- 3. Central limit theorem**

# Central Limit Theorem (CLT)

Central limit theorem (CLT):

If you have a population with mean  $\mu$  and take **sufficiently large random samples** from the population **with replacement**, then the distribution of **the sample means** will be **approximately normally distributed**.

- Useful because we **never know** the population's shape!
- Now if there is a large enough sample, we can **model the sampling distribution for the mean**.

If you have a population with mean  $\mu$  and take **sufficiently large random samples** from the population **with replacement**, then the distribution of **the sample means** will be **approximately normally distributed**.

也就是说如果我们的样本满足两个条件

一个是 sample size 足够大

另一个是，抽样的方式是有放回的

那么  $x\bar{}$  的 sampling Distribution 就是可以近似看作 Normal Distribution

这个定理非常有用，因为大多数情况我们都不知道 population 的分布是什么样的。那在这种情况下，我依旧可以知道 sample mean 的分布大概什么样子，只要满足了刚刚提到的两个条件，我们就可以把 mean 的 sampling Distribution 当做 Normal 来处理。这个定理也是我们后面做假设检验的一个最基础的理论依据。非常有意义。

## Condition 1: Random

You must ensure the sample was **randomly selected** from the **target population** (in the case of an experiment, you must ensure there was a random assignment to treatment)

### Why must this condition be satisfied?

If this condition is not satisfied, your estimator is **biased**.

If it is satisfied, your sampling distribution is **centered** at the true population mean value →  $\mu_{\bar{x}} = \mu$

我们昨天得出了 $\bar{x}$ 的分布，但是有一些条件，我们今天来具体的看一看  
第一个就是 **random sampling method**

这个条件满足的话，我们就会得到一个**unbiased sample**，那也就得到了**unbiased estimator**

所以就可以保证  $E(\bar{X}) = \mu$

## Condition 2: 10% Condition

### a) What is the condition?

The **sample size** ( $n$ )

**size (N)**:  $n < 0.10N$

If sampling without replacement, the standard deviation of the sample mean is smaller than what is given by the formula above. If the sample size is less than 10% of the population size, the difference is negligible.

### b) Why must this condition be satisfied?

It ensures this is true:  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

For example if sample = population,  $\sigma_{\bar{x}} = ?$

第二个条件是用来保证variance的，上节课我们说到，如果抽样的时候是with replacement的话， $\bar{x}$ 的方差是可以求出来的，但是如果放回抽样的话，就比较复杂，可能很难算出方差是多少，但是现实中抽样基本都是without replacement，但如果满足这个10% condition的话，也就是population size 比sample size大很多的话，就可以近似成with replacement。

我们也可以举个反例，如果sample很大，大到趋近于population的话， $\bar{x}$ 就变成了mu，那是不是就失去了波动性？

## Condition 2: 10% Condition

### a) What is the condition?

The **sample size (n)** must be **less than 10% of the population size (N)**:  $n < 0.10(N)$

### b) Why must this condition be satisfied?

It ensures this is true:  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

For example if sample = population,  $\sigma_{\bar{x}} = 0$

## Condition 3: Normal/Large Sample

### a) What is the condition?

- i) The **population** distribution must be **normal OR**
- ii) **The sample size is 30 or more** ( $n \geq 30$ )

### b) Why must this condition be satisfied?

This provides evidence that the sampling distribution is **approximately normal** in shape.

For example if sample size = 1, sampling distribution = ?

第三个也就是最后一个条件，是来保证 $\bar{x}$ 的分布是正态分布的  
这里面有两个条件都是可以的

第一个，如果population Distribution是normal的话，那无论sample size是多少，  
我们得到的sample mean都是正态分布的（可以详细解释一下）

如果我们不知道population Distribution或者它不是正态分布的话，就需要满足  
sample size larger than or equal to 30这个条件，也就是用到了central limit theory

我们可以举个反例说明一下

假设现在的sample size=1，那 sampling Distribution是啥？

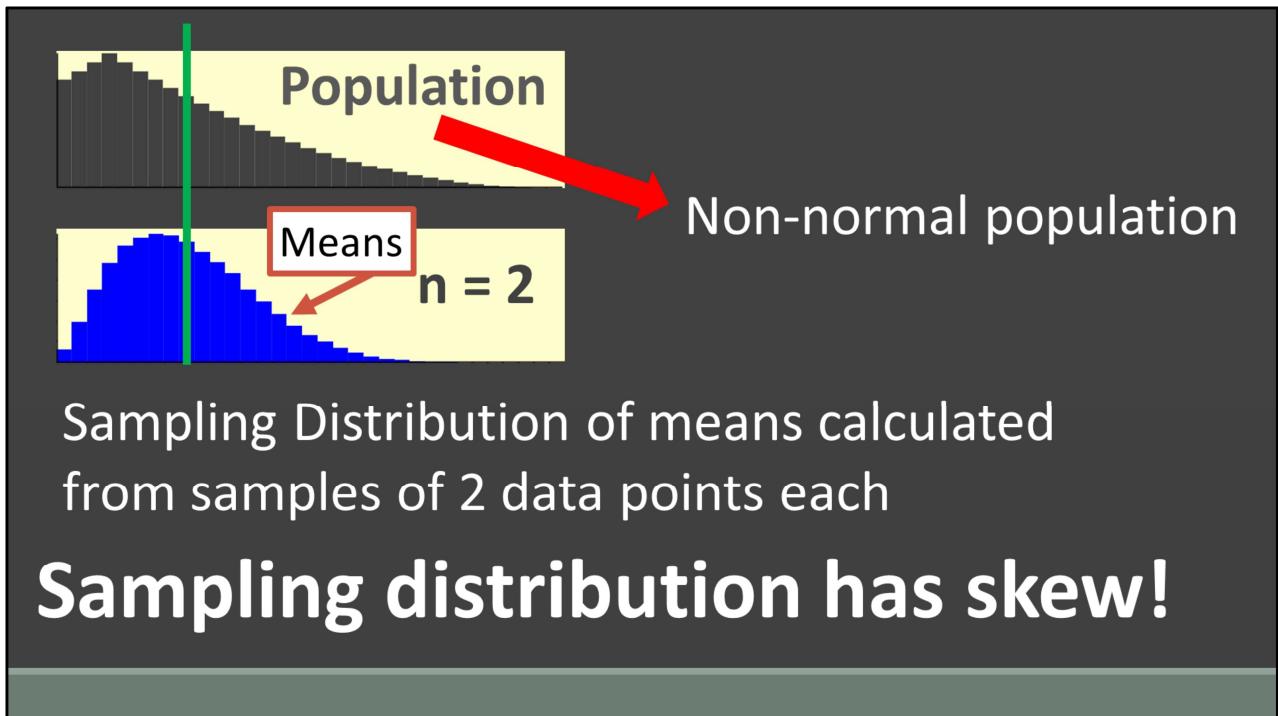
## Condition 3: Normal/Large Sample

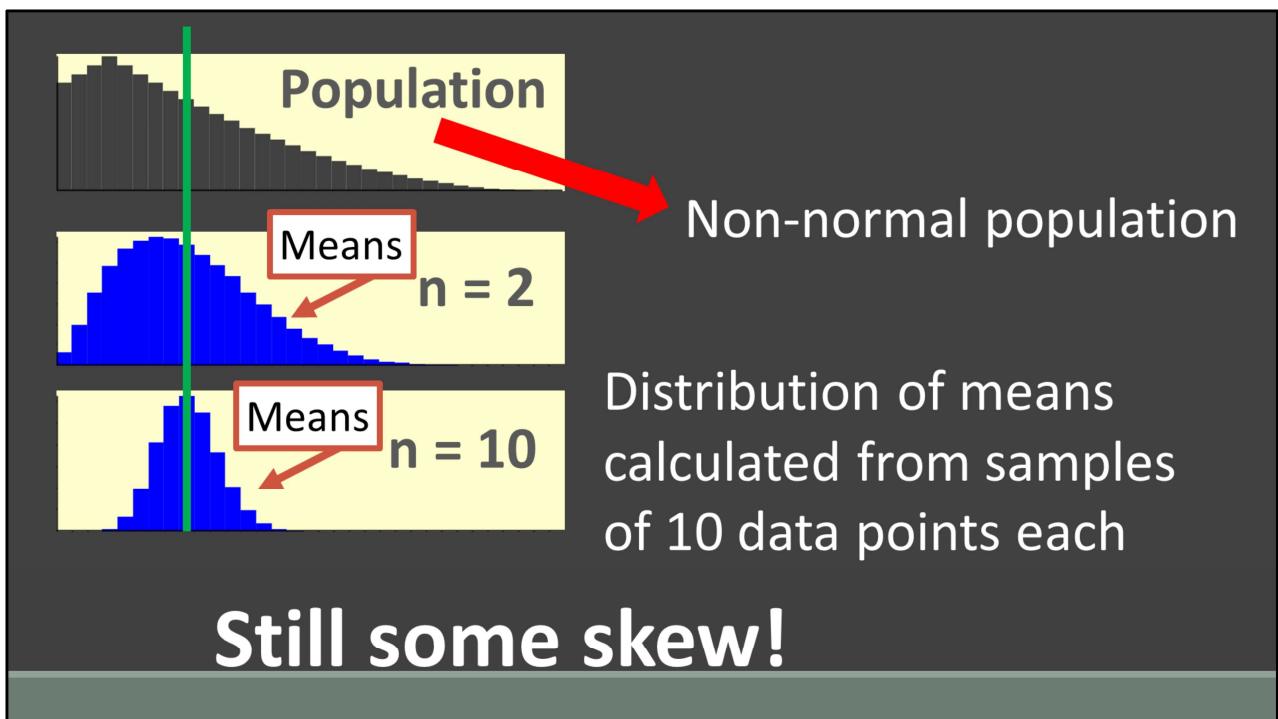
### a) What is the condition?

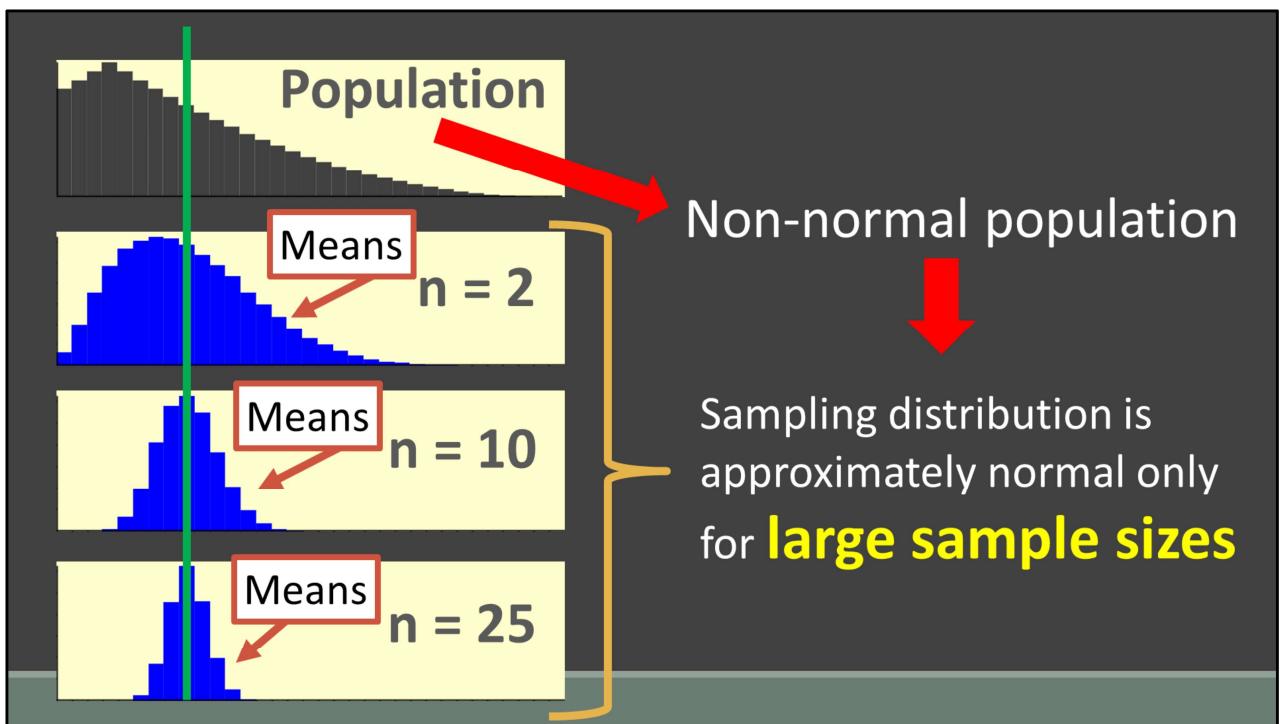
- i) The **population** distribution must be **normal OR**
- ii) The sample size is **30 or more** ( $n \geq 30$ )

### b) Why must this condition be satisfied?

For example if sample size = 1, sampling distribution = **population distribution.**







所以Sampling Distribution是一个随着sample size的变大，逐渐趋近于Normal

# Summary

$$\bar{X} \sim Norm \left( \mu_{\bar{X}} = \mu, \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \right)$$

3) Normal/Large Sample  
→ approx. normal  
**shape** (by CLT)

1) Random condition  
→ unbiased **center**

2) 10% condition  
→ calculable **spread**

我们来总结一下

对于 $\bar{X}$ 的分布，我们需要满足几个条件

第一个，如果保证它的期望就是 $\mu$ 的话，需要满足random sample

如果要保证它的方差可以算出来的话，需要满足with replacement或者10% condition，一般来说都是无法满足without replacement，但是这里考点比较灵活，大家也最好知道一下原理，原则上如果with replacement的话，我们的方差是可以保证能算出来的，10% condition只是一个近似。

第三个条件是针对normal的，如果population就是normal，那 $\bar{X}$ 就是完美的normal

如果不是，但是满足了 $n \geq 30$ ，那 $\bar{X}$ 就可以趋近于normal

这三个条件分别对应了三部分，期望，方差和分布类型

这三个条件是相互独立的，比如，如果第一个random sample的条件满足了，其他的都没满足，我还是可以得到 $\bar{X}$ 的期望是 $\mu$ 这个结论的。

# Sampling distributions for differences in sample means

For a numerical variable, when randomly sampling with replacement from two independent populations with population mean  $\mu_1$  and  $\mu_2$  and population standard deviation  $\sigma_1$  and  $\sigma_2$ , the sampling distribution of the difference in sample means  $\bar{X}_1 - \bar{X}_2$  has mean  $\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2$  and standard deviation

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

The sampling distribution of the difference in sample means  
 $\overline{X_1} - \overline{X_2}$  can be modeled with a normal distribution **IF**  
the two sampling distributions can be modeled with a normal  
distribution.

If  $X \sim N(\mu_X, \sigma_X^2)$ ,  $Y \sim N(\mu_Y, \sigma_Y^2)$ ,  
X and Y are independent random variables,  
then  $X+Y$  follows the normal distribution.

The sampling distribution of the difference in sample means

$\bar{X}_1 - \bar{X}_2$  can be modeled with a normal distribution **IF**

the two sampling distributions can be modeled with a normal distribution:

1. if the population distributions can be modeled with a normal distribution.
2. if cannot, sample sizes have to be greater than or equal to 30.