

## Sampling Distribution and the Central Limit Theorem

Goal: study the distribution of statistics obtained using a random sample

Let  $\{x_1, x_2, \dots, x_n\}$  be a random sample of size  $n$ .  
 (1)  $x_i$ 's have the same distribution for  $i=1, 2, \dots, n$

(2) The set  $\{x_1, x_2, \dots, x_n\}$  is an independent set of random variables.

defn:  $X = \text{value of the random variables in the random sample, call it } X = \{x_1, x_2, x_3, \dots, x_n\}$

The joint sample space for the r.v.s  $x_1, x_2, \dots, x_n$  is the Cartesian product of  $X$ .

$$X^n = \{(x_1, x_2, \dots, x_n); x_i \in \mathbb{R}\}$$

Since  $x_1, x_2, \dots, x_n$  are independent r.v.s

The joint probability distribution of  $x_1, x_2, x_3, \dots, x_n$

(1)  $p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$   
 if  $x_i$ 's are discrete, and  
 $p(x_i)$  is the common prob of all  $x_i$ 's.

(2)  $f(x_1, x_2, \dots, x_n) = f(x_1)f(x_2) \cdots f(x_n)$   
 if  $x_i$ 's are continuous and  
 $f(x_i)$  is the common pdf of all  $x_i$ 's.

For brevity, we will work with discrete r.v.s.  
 Continuous case is similar.

Let  $T$  be a statistic for the random sample  $\{x_1, x_2, \dots, x_n\}$   
 random variable calculated using r.v.s  $x_1, x_2, \dots, x_n$

$T$  is a function defined

$T: X^n \rightarrow \mathbb{R}$ , where  
 $X^n = \text{joint sample space of } x_1, x_2, \dots, x_n$

The probability distribution of the statistic  $T$   
 ↓  
 Sampling distribution of  $T$

To calculate the sampling distribution of  $T$   
 → use the joint distribution of the r.v.s  $x_1, x_2, \dots, x_n$

Let  $p(x_1, x_2, \dots, x_n)$  be the joint prob of  $x_1, x_2, \dots, x_n$   
 (we are assuming the  $x_i$ 's are discrete)

Want: the b.m.f. of  $T$  → i.e. for  $t \in \mathbb{C}$ , want  
 $P(T=t)$

Let  $E(\cdot) = \{x_1, x_2, \dots, x_n\} \subset \mathbb{C}^n$ ,  $T(x_1, x_2, \dots, x_n) = t$ ,

then

$$\begin{aligned} P(T=t) &= \sum_{(x_1, x_2, \dots, x_n) \in E(\cdot)} p(x_1, x_2, \dots, x_n) \\ &= \sum_{(x_1, x_2, \dots, x_n) \in E(\cdot)} p(x_1)p(x_2) \cdots p(x_n) \end{aligned}$$

where  $p(x)$  is the common prob of  $x_1, x_2, \dots, x_n$ .

The Sample mean and sample total

Let  $X = \{x_1, x_2, \dots, x_n\}$  → random sample coming from a population with mean =  $\mu$  variance =  $\sigma^2$ .

ie  $E(X_i) = \mu$  for  $i=1, 2, \dots, n$   
 $V(X_i) = \sigma^2$  for  $i=1, 2, \dots, n$

The "sample total" for  $\{x_1, x_2, \dots, x_n\}$  →  $T_0 = x_1 + x_2 + \dots + x_n$

The "sample mean" for  $\{x_1, x_2, \dots, x_n\}$  →  $\bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n}$

Can calculate:

$$\begin{aligned} E(T_0) &= E(x_1 + x_2 + \dots + x_n) \\ &= E(x_1) + E(x_2) + \dots + E(x_n) \\ &\quad \vdots \quad \vdots \\ &= n\mu \end{aligned}$$

$$\text{ie } E(T_0) = n\mu$$

Similarly,

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

Question: What is the sampling distribution of  $\bar{X}$  and  $T_0$ ?  
 ↓  
 ↓

The central limit theorem provides the answer!

Central Limit Theorem (Normal Case)

Suppose  $\{x_1, x_2, \dots, x_n\}$  is a random sample coming from a population with normal distribution  $N(\mu, \sigma^2)$ , then

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

$$T_0 \sim N(n\mu, n\sigma^2)$$

Equivalently, the random variables

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \quad \text{and} \quad \frac{T_0 - n\mu}{\sigma\sqrt{n}}$$

have the standard normal distribution.

Note: For CLT to hold for all  $n \in \mathbb{N}$  (especially when  $n$  is small)  
 ↓ need population to be normally dist.

Can drop normality assumption if  $n$  is large.

Central Limit Theorem (General Case)

If  $\{x_1, x_2, \dots, x_n\}$  is a random sample from a population with mean =  $\mu$  variance =  $\sigma^2$ , then

if  $n$  is large enough, the distribution of  $\bar{X}$  and  $T_0$  is approximately normal.

$$\text{i.e. } \bar{X} \sim N(\mu, \frac{\sigma^2}{n}), T_0 \sim N(n\mu, n\sigma^2)$$

The approximation improves as  $n \rightarrow \infty$ .

The distribution of  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  and  $\frac{T_0 - n\mu}{\sigma\sqrt{n}}$  approaches the standard normal distribution as  $n \rightarrow \infty$ .

Example: Suppose the population =  $\{0, 1\}$

The distribution of this population is the Bernoulli dist.

$$\begin{aligned} P(X=1) &= P(\text{success}) = p = P(1) \\ P(X=0) &= P(\text{failure}) = 1-p = P(0) \end{aligned}$$

Let  $\{x_1, x_2, \dots, x_n\}$  be a random sample from  $\{0, 1\}$

Goal: Calculate the sampling distribution of  $T_0$

Suppose  $n=2$ , then  $\{x\} = \{(0, 0), (1, 0), (0, 1), (1, 1)\}$

Value of  $T_0 = \{0, 1, 2, 3\}$

$p(x,y) = p(x)p(y)$  → be the joint prob of  $x_1$  and  $x_2$

then  $P(T_0=0) = p(0,0) = p(0)p(0) = (1-p)^2$

$P(T_0=1) = p(0,1) + p(1,0) = p(0)p(1) + p(1)p(0) = 2p(1-p)$

$P(T_0=2) = p(1,1) = p(1)p(1) = p^2$

Putting this into a table gives

$T_0$	0	1	2
$P(T_0=t)$	$(1-p)^2$	$2p(1-p)$	$p^2$

notice anything strange?

Suppose  $n=3$ :

Now  $\{x\} = \{(0,0,0), (1,0,0), (0,1,0), (1,1,0), (0,0,1), (1,0,1), (0,1,1), (1,1,1)\}$

$p(x_1, x_2, x_3) = p(x_1)p(x_2)p(x_3)$  and

Value of  $T_0 = \{0, 1, 2, 3\}$

$P(T_0=0) = p(0,0,0) = p(0)p(0)p(0) = (1-p)^3$

$P(T_0=1) = p(0,0,1) + p(0,1,0) + p(1,0,0) = p(0)p(0)p(1) + p(0)p(1)p(0) + p(1)p(0)p(0) = 3p^2(1-p)$

$P(T_0=2) = p(0,1,1) + p(1,0,1) + p(1,1,0) = p(0)p(1)p(1) + p(1)p(0)p(1) + p(1)p(1)p(0) = 3p^2(1-p)^2$

$P(T_0=3) = p(1,1,1) = p(1)p(1)p(1) = p^3$

The probability distribution table of  $T_0$  is

$T_0$	0	1	2	3
$P(T_0=t)$	$(1-p)^3$	$3p^2(1-p)$	$3p^2(1-p)^2$	$p^3$

Suppose  $n=4$ :

Now  $\{x\} = \{(0,0,0,0), (1,0,0,0), (0,1,0,0), (1,1,0,0), (0,0,1,0), (1,0,1,0), (0,1,1,0), (1,1,1,0), (0,0,0,1), (1,0,0,1), (0,1,0,1), (1,1,0,1), (0,0,1,1), (1,0,1,1), (0,1,1,1), (1,1,1,1)\}$

$(0,0,0,0), (0,0,1,1), (0,1,0,1)$

The joint dist is  $p(x_1, x_2, x_3, x_4) = p(x_1)p(x_2)p(x_3)p(x_4)$

Value of  $T_0 = \{0, 1, 2, 3, 4\}$

The distribution of  $T_0$  can be calculated as

$T_0$	0	1	2	3	4
$P(T_0=t)$	$(1-p)^4$	$4(1-p)^3p$	$6(1-p)^2p^2$	$4p^3(1-p)$	$p^4$

In general for arbitrary  $n$

Value of  $T_0 = \{0, 1, 2, \dots, n\}$  and

$P(T_0=k) = \binom{n}{k} (1-p)^{n-k} p^k$   $k=0, 1, 2, \dots, n$

↓  
 prob of the Binomial distribution with parameters  $n$  and  $p$

↓  
 $T_0 \sim \text{Bin}(n, p)$  → makes sense because  
 ↑ To counts the # of '1's (success)  
 for given sample of size  $n$ , where  
 each sample entry is an outcome of  
 independent Bernoulli trial

The Binomial Experiment

↓  
 $T_0 \sim \text{Bin}(n, p) \Rightarrow E(T_0) = np$   
 $V(T_0) = np(1-p)$

Also the central limit theorem says,  
 for large enough  $n$  →  
 $T_0 \sim N(n\mu, n\sigma^2)$   
 $N(np, np(1-p))$

↓  
 where  
 $\mu = \text{mean}(X) = p$   
 $\sigma^2 = \text{Var}(X) = p(1-p)$

↓  
 $T_0 \sim \text{Bin}(n, p) \Rightarrow$  for large  $n$   
 $\text{Bin}(n, p) \sim N(np, np(1-p))$

↓  
 Normal approximation for Binomial!