

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 $E(\cdot) = \{x_1, x_2, \dots, x_n\} \in X^n$, $T(x_1, x_2, \dots, x_n) = t$,

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 = μ variance = σ^2 .

i.e. $E(X_i) = \mu$ 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:

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

$$\text{i.e. } E(T_0) = n\mu$$

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 $N(\mu, \sigma^2)$

for any $n \in \mathbb{N}$ the distributions of the sample mean \bar{X} and the sample total T_0 are

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

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

Goal: Calculate the sampling distribution of T_0

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

Values of $T_0 = \{0, 1, 2\}$

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

then $P(T_0=0) = p(0,0) = p_1(0)p_2(0) = (1-p)^2$

$P(T_0=1) = p(0,1) + p(1,0) = p_1(0)p_2(1) + p_1(1)p_2(0) = 2p_1(1)p_2(0)$

$P(T_0=2) = p(1,1) = p_1(1)p_2(1) = p^2$

Putting this into a table gives

T_0	0	1	2
$P(T_0=t)$	$(1-p)^2$	$2p_1(1)p_2(0)$	p^2

Notice anything strange?

Suppose $n=3$:

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

$p(x_1, x_2, x_3) = p_1(x_1)p_2(x_2)p_3(x_3)$ and

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

$P(T_0=0) = p(0,0,0) = p_1(0)p_2(0)p_3(0) = (1-p)^3$

$P(T_0=1) = p(1,0,0) + p(0,1,0) + p(0,0,1) = p_1(1)p_2(0)p_3(0) + p_1(0)p_2(1)p_3(0) + p_1(0)p_2(0)p_3(1) = 3p_1(1)p_2(0)p_3(0)$

$P(T_0=2) = p(1,1,0) + p(1,0,1) + p(0,1,1) = p_1(1)p_2(1)p_3(0) + p_1(1)p_2(0)p_3(1) + p_1(0)p_2(1)p_3(1) = 3p_1(1)p_2(1)p_3(0)$

$P(T_0=3) = p(1,1,1) = p_1(1)p_2(1)p_3(1) = p^3$

The probability distribution table of T_0 is

T_0	0	1	2	3
$P(T_0=t)$	$(1-p)^3$	$3p_1(1)p_2(0)p_3(0)$	$3p_1(1)p_2(1)p_3(0)$	p^3

Suppose $n=4$:

Now $X = \{(0,0,0,0), (1,0,0,0), (1,1,0,0), (0,1,0,0), (0,0,1,0), (1,1,1,0), \dots, (0,0,0,1), (0,0,1,1), (0,1,1,1)\}$

\vdots

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

The joint dist of $p(x_1, x_2, x_3, x_4) = p_1(x_1)p_2(x_2)p_3(x_3)p_4(x_4)$

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

The distribution of T_0 can be calculated as

T_0	0	1	2	3	4
$P(T_0=t)$	$(1-p)^4$	$4(1-p)^3p$	$6(1-p)^2p^2$	$4(1-p)p^3$	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
 $\mu = \text{mean}(X) = p$
 $\sigma^2 = \text{Var}(X) = p(1-p)$

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

Normal approximation for Binomial!

$\{X_1, X_2, X_3, \dots, X_n\}$ \rightsquigarrow (i) X_i 's have the same dist
 $\underbrace{\qquad}_{\text{random sample of size } n}$ (ii) X_1, X_2, \dots, X_n are ind r.v.'s.

A statistic \rightsquigarrow random variable calculated using a random sample.

$\underbrace{\qquad}_{\text{can calculate the distribution of the statistic}}$
 $\underbrace{\qquad}_{\text{"sampling distribution".}}$

Suppose $X = \text{set of values of r.v. in } \{X_1, \dots, X_n\}$

joint sample space = $X^n = \{(x_1, x_2, \dots, x_n) \mid x_i \in X\}$
 of X_1, X_2, \dots, X_n

can calculate the "joint probability distribution"

Assume that the X_i 's are discrete \rightsquigarrow want the joint pmf of X_1, X_2, \dots, X_n .

Joint pmf: $p: X^n \rightarrow [0, 1]$

$$p(x_1, x_2, x_3, \dots, x_n) = p_{x_1}(x_1) \cdot p_{x_2}(x_2) \cdots p_{x_n}(x_n)$$

where $p_{x_i}(x_i) \rightsquigarrow$ pmf of the common dist of all X_i 's

$$p: X \rightarrow [0, 1].$$

Let T be a statistic $\rightsquigarrow T: X^n \rightarrow \mathbb{R}$

To calculate the pmf $T \rightsquigarrow$ want $P(T=t)$

for $t \in \mathcal{T} = \text{values of } T$.

(i) Find out what the values of T are.

(ii) Use the joint pmf to calculate $\underline{P(T=t)}$

$$t \in \mathcal{T}, E(t) = \{(x_1, x_2, \dots, x_n) \in X^n \mid T(x_1, x_2, \dots, x_n) = t\}$$

$$= T^{-1}(t) \subseteq X^n$$

Since $E(t) \subseteq X^n \rightsquigarrow$ know how to calculate $P(E(t))$

$$\therefore P(T=t) = P(E(t)) = \sum_{(x_1, x_2, \dots, x_n) \in E(t)} p(x_1, x_2, \dots, x_n).$$

Example: Population = $\{0, 1\}$

$$p_x(0) = p, p_x(1) = 1-p$$

x	0	1
$p_{x(0)}$	$(1-p)$	p
$p_{x(1)}$		

Get a random sample of size 2

$\{X_1, X_2\}$. X_1 and X_2 have the Bernoulli dist defined x .

$$X = \{0, 1\}$$

Joint sample space = $\{(0,0), (0,1), (1,0), (1,1)\}$

Joint pmf: $p(x,y) = p_x(x) \cdot p_y(y)$

$$\Rightarrow p(0,0) = p_x(0) \cdot p_y(0) = (1-p)(1-p) = (1-p)^2$$

$$p(0,1) = p_x(0) \cdot p_y(1) = (1-p) \cdot p$$

$$p(1,0) = p_x(1) \cdot p_y(0) = p \cdot (1-p)$$

$$p(1,1) = p_x(1) \cdot p_y(1) = p^2$$

x	y	0	1
0	0	$(1-p)^2$	$p(1-p)$
1	1	$p(1-p)$	p^2

$$\text{Sample mean} := \bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{T_0}{n}$$

$$\text{Sample Total} := T_0 = x_1 + x_2 + \dots + x_n$$

$$\mathcal{T} = \{(0,0), (0,1), (1,0), (1,1)\}$$

Want: values of \bar{X} and T_0

$$\text{Values of } \bar{X} = \{0, \frac{1}{2}, 1\}$$

$$X(0,0) = \frac{0+0}{2} = 0, X(1,0) = X(0,1) = \frac{1}{2}, X(1,1) = \frac{1+1}{2} = 1$$

$$T_0(0,0) = 0, T_0(0,1) = T_0(1,0) = 1+0 = 1, T_0(1,1) = 1+1 = 2$$

$$\text{Values of } T_0 = \{0, 1, 2\}$$

To calculate the distribution of \bar{X} and T_0

use the joint pmf $p(x_1, x_2)$

$$E(0) = \{0, 0\}$$

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

$$E(1) = \{0, 1, 1, 0\}$$

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

$$E(2) = \{1, 1\}$$

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

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

the dist table for
 $\text{Bin}(n=2, p=p)$

Calculating the distribution table for \bar{X} .

$$\text{Values of } \bar{X} = \{0, \frac{1}{2}, 1\}$$

$$P(\bar{X}=0) = p(0,0) = (1-p)^2$$

$$P(\bar{X}=\frac{1}{2}) = P(T_0=1) = 2(1-p)p$$

$$P(\bar{X}=1) = P(T_0=2) = p^2$$

\bar{x}	0	$\frac{1}{2}$	1
$P(\bar{X}=\bar{x})$	$(1-p)^2$	$2p(1-p)$	p^2

General formula for the sampling dist of \bar{X} and T_0

using the Central Limit theorem.

I) CLT (Normal case)

If $\{x_1, x_2, \dots, x_n\}$ coming from a normal dist

$$N(\mu, \sigma^2)$$

for all $n \in \mathbb{N}$

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

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

can show that

$$E(\bar{X}) = \mu \quad V(\bar{X}) = \frac{\sigma^2}{n} \Rightarrow \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

$$E(T_0) = n\mu \quad V(\bar{X}) = n\sigma^2 \Rightarrow \sigma_{T_0} = \sqrt{n} \sigma$$

CLT in normal case $\Rightarrow \frac{\bar{X}-\mu}{\sigma/\sqrt{n}}$ and

$$\frac{T_0 - n\mu}{\sqrt{n}\sigma}$$

both have the standard normal dist.

II) CLT (General Case)

Suppose $\{x_1, x_2, \dots, x_n\}$ is

a random sample coming from a population with

$$\text{Mean} = \mu$$

$$\text{Variance} = \sigma^2$$

If n is large
 \bar{X} and T_0 have approximately normal distributions.

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

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

as $n \rightarrow \infty$, the approx gets better

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

standard normal dist as $n \rightarrow \infty$.

Example (Revisited).

Suppose $n=8$ of size n coming $\{0, 1\}$

$$\{x_1, x_2, x_3, \dots, x_n\} \quad x_i \in \{0, 1\}$$

$$\text{sample data} = X^n = \{x_1, x_2, \dots, x_n\} \mid x_i \in \{0, 1\}\}$$

$$|X^n| = 2^2 \times 2^2 \times 2^2 \times 2^2 = 2^8$$

typical element in $X^n \rightarrow (0, 0, 1, 1, 0, 0, \dots)$

values of $T_0 : \{0, 1, 2, 3, \dots, n\}$

$$T_0(1, 1, \dots, 1) = 1+1+1+\dots+1 = n$$

$$T_0(0, 0, \dots, 0) = 0+0+\dots+0 = 0$$

$$\{T_0 = k\} = \{(x_1, \dots, x_n) \mid k \text{ } x_i \text{'s are equal}\}$$

to 1

$$\Rightarrow |\{T_0 = k\}| = \binom{n}{k}$$

$$P(\underbrace{1, 1, 1, \dots, 1}_{k}, 0, 0, 0) = p^k \cdot (1-p)^{n-k}$$

$$\therefore P(T_0 = k) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$$

\downarrow
pmf of $\text{Bin}(n, p)$

$\therefore T_0 \sim \text{Bin}(n, p) + \text{CLT} \Rightarrow \text{for large } n$

$$\text{Bin}(n, p) \sim N(np, np(1-p))$$

=====

Values of \bar{X} for $\{x_1, x_2, \dots, x_n\}$

Values of T_0 : $\{0, 1, 2, \dots, n\}$

Values of \bar{X} : $\{0, \frac{1}{n}, \frac{2}{n}, \frac{3}{n}, \dots, 1\}$

