

Random Sample and Statistics

Collecting data → important aspect of doing statistics.

we want to use a subset of the population to draw conclusion/inference about the entire population.

this subset is called a "sample".

the way a sample is obtained will affect our confidence in the conclusions we derive using it.

Example: Suppose we have a vaccine for Covid-19

we would be interested in knowing its "efficacy".

Consider the following experiment:

- ① Get two samples of people, group A and group B.
- ② Give vaccine to group A → treatment arm
placebo to group B → control arm
- ③ Expose both groups to the virus, and measure the appropriate response variable.
- ④ Apply statistical tests to check if there is statistically significant difference between the two responses.

To get correct conclusions about the vaccine → group A and B need to be carefully selected

(i) Is it possible that certain people can never make it to the samples?

(ii) Do your sample contain only young individuals?

(iii) Are your samples representative of the population?

(iv) Are the sample sizes large enough?

Statistical tools → usually based on principles of probability

need a precise notion of "randomness" of a sample.

Random Sample

We say $\{X_1, X_2, \dots, X_n\}$ a set of random variables is a "random sample"

if:

(i) X_i has the same probability distribution for $i=1, 2, \dots, n$.

(ii) the set $\{X_1, X_2, \dots, X_n\}$ is an independent set of random variables.

Note ① Sometimes a random sample

also referred to as a set of independent and identically distributed random variables.

"i.i.d" random variables.

② The common distribution of all the X_i 's

"population distribution".

③ $\{X_1, X_2, X_3, \dots, X_n\}$ → sampling procedure

Random sample
a collection of random variables

Sample data
Changes everytime we run the sampling procedure.

Example: Suppose we are getting a random sample of size 3

from the population: $\{1, 2, 3, 4, 5\}$

To get a random sample → need to sample with replacement.

Let $\{Num1, Num2, Num3\}$ → random sample

Sample first number in the sample
Sample 2nd number in the sample
Collection of r.v's
Sample 3rd number in the sample

∴ Sampling with replacement → there are $5 \times 5 \times 5$ ways of getting sample data.

ie. we have S^3 choices of sample datasets

$$\{(1,1,1), (1,1,2), (1,2,1), (2,1,1), \dots, (5,5,5)\}$$

Each run of sampling procedure gets one of the elements in

Statistic

We are often interested in some features/characteristics of the random sample

A statistic → is a quantity calculated using a random sample (say $\{X_1, X_2, \dots, X_n\}$)

Since the calculation of the statistic involves random variables

the statistic will also be a random variable.

Note: ① The value of the statistic changes as the sample data changes.

② A statistic is a random variable defined on sample data corresponding to a random sample.

If $X = \text{Value of the r.v's in random sample}$

Sample data → $\{(x_1, x_2, \dots, x_n) | x_i \in \mathbb{R}, i=1, \dots, n\}$

Since $\{X_1, X_2, \dots, X_n\}$ is an independent set, the "joint distribution" for X_1, X_2, \dots, X_n is given as

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

where $p_x(x)$ → pmf of the common distribution X .

We are assuming that X is discrete.

If X is continuous use pdf $f_X(x)$ to define the joint pdf of X_1, X_2, \dots, X_n .

③ Can use the joint pmf/pdf of X_1, X_2, \dots, X_n to calculate the probability distribution of the statistic of interest

"Sampling distribution" of the statistic.

Example Revisited:

Recall: Population = $\{1, 2, 3, 4, 5\}$

Random sample = $\{Num1, Num2, Num3\}$

Examples of statistics defined using $\{Num1, Num2, Num3\}$

Name of the statistic	Definition of the statistic	Example of evaluation on sample data
-----------------------	-----------------------------	--------------------------------------

① Sample mean	$\bar{X} = \frac{Num1 + Num2 + Num3}{3}$	$\bar{X}(1, 1, 1) = \frac{1+1+1}{3} = 1$ $\bar{X}(2, 1, 2) = \frac{2+1+2}{3} = \frac{5}{3}$ $\bar{X}(3, 1, 1) = \frac{3+1+1}{3} = \frac{5}{3}$ $\bar{X}(4, 5, 1) = \frac{4+5+1}{3} = \frac{10}{3}$
---------------	--	---

② Sample Median	$\tilde{X} = \text{median}\{\text{Num1, Num2, Num3}\}$	$\tilde{X}(1, 1, 1) = 1$ $\tilde{X}(2, 1, 2) = 2$ $\tilde{X}(3, 3, 1) = 3$ $\tilde{X}(1, 5, 4) = 4$
-----------------	--	--

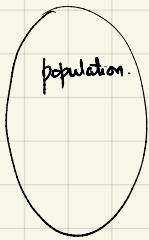
③ Sample Maximum	$\text{Max} = \max\{\text{Num1, Num2, Num3}\}$	$\text{Max}(1, 1, 1) = 1$ $\text{Max}(1, 2, 5) = 5$ $\text{Max}(4, 1, 3) = 4$ $\text{Max}(3, 1, 2) = 3$
------------------	--	--

④ Sample Total	$T_0 = \text{Num1} + \text{Num2} + \text{Num3}$	$T_0(1, 1, 1) = 1+1+1 = 3$ $T_0(1, 3, 5) = 1+3+5 = 9$ $T_0(2, 2, 4) = 2+2+4 = 8$ $T_0(3, 1, 2) = 3+1+2 = 6$
----------------	---	--

⑤ Sample Variance	$S^2 = \frac{(Num1 - \bar{X})^2 + (Num2 - \bar{X})^2 + (Num3 - \bar{X})^2}{2}$	$S^2(1, 1, 1) = \frac{(1-\bar{X}(1, 1, 1))^2 + (1-\bar{X}(1, 1, 1))^2 + (1-\bar{X}(1, 1, 1))^2}{2} = 0$ $S^2(1, 2, 3) = \frac{(1-\bar{X}(1, 2, 3))^2 + (2-\bar{X}(1, 2, 3))^2 + (3-\bar{X}(1, 2, 3))^2}{2} = \frac{(1-2)^2 + (2-2)^2 + (3-2)^2}{2} = \frac{-1 + 0 + 1}{2} = 0$
-------------------	--	---

To calculate the sampling distribution of these statistics
Need to know the distribution of the population $\{1, 2, 3, 4, 5\}$

Random Samples and Statistics



sample-data.

$$\{x_1, x_2, x_3, \dots, x_n\}$$

{Student 1, Student 2, Student 3, Student 4, Student 5}

measure their heights.

Sample 1: {John, Macelius, Mia, Neo, Axmin}

Sample data: {5.9, 5.12, 5.6, 4.2, 6.0}

Sample 2: {Kevin, Nina, Rhea, Sounath, Angula}

Sample 3: {Matt, Sofia, Luke, Smith, Sergeant}

Sampling without replacement vs Sampling with replacement.



every subsequent sample choice depends on the other choices

choices are going to be independent of each other.

Sample 4: {Dennis, Dennis, Dennis, Dennis, Dennis}

Typically: If the sample is small relative to the population size



we can say that "Sampling without replacement will be approximately the same as sampling with replacement."



If sample is ≤ 0.05 of the population.

Sampling with replacement: gives sample choices that are independent of each other. (good)

might not get a sample of "distinct" elements.

(not so good!)

test subjects

Sampling without replacement: gives sample choices that are dependent on each other (bad)

gives a sample of distinct individuals.

Random Sample \rightarrow collection of n random variables
size n

$$\{x_1, x_2, x_3, \dots, x_n\}$$

(i) The probability dist of X_i for $i=1, 2, \dots, n$ is the same.

distribution of the population.

(ii) $\{x_1, x_2, \dots, x_n\}$ is a independent set random variable.

Recall: $\{x_1, x_2, \dots, x_n\}$ are ind r.v.s if



$$\{x_1, x_2, \dots, x_n\} \subseteq \{x_1, \dots, x_n\}$$

the joint probability dist of $x_{i_1}, x_{i_2}, \dots, x_{i_n}$ satisfies

$$P(x_1, x_2, \dots, x_n) = P_{x_1}(x_1) \cdot P_{x_2}(x_2) \cdots P_{x_n}(x_n)$$

Random sample $\{x_1, x_2, x_3, \dots, x_n\}$ $\xrightarrow{\text{Sampling procedure}}$ $(x_1, x_2, x_3, \dots, x_n)$
random variables \downarrow sample data

Example:

Population: {2, 5, 9, 10}

Want a random sample of size 2. (need to sample with replacement)

Sample = {Num1, Num2} $\xrightarrow{\text{Sampling procedure}}$ sample data

$$\{(x_1, x_2) | x_1, x_2 \in \{2, 5, 9, 10\}\}$$

\therefore Choices of sample $\xrightarrow{\text{Sampling procedure}}$ data

$4^2 = 16$ choices of sample data

$$\{(2,2), (2,5), (5,2), (2,9), \dots, (10,10)\}$$

Statistic \rightarrow Some number calculated
using a random sample.

\Downarrow
Statistic itself will be a
random variable.
 \Downarrow
Value of the statistic changes
everytime the sample data
changes.

Num1\Num2	2	5	9	10
2	(2,2)	(2,5)	(2,9)	(2,10)
5	(5,2)	(5,5)	(5,9)	(5,10)
9	(9,2)	(9,5)	(9,9)	(9,10)
10	(10,2)	(10,5)	(10,9)	(10,10)

Example: Population: $\{2, 5, 9, 10\}$ —

$x =$	2	5	9	10
$P(x)$	0.1	0.1	0.4	0.2

Sample: $\{\text{Num1}, \text{Num2}\}$

Interesting Statistics

① Sample mean.

$$\bar{X} = \frac{\text{Num1} + \text{Num2}}{2}$$

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

$$\bar{X}(5,9) = \frac{5+9}{2} = \frac{14}{2} = 7$$

$$\bar{X}(5,10) = \frac{5+10}{2} = \frac{15}{2} = 7.5$$

evaluating the statistic
on sample data.

② Sample Min

$$\text{Min} = \min \{\text{Num1}, \text{Num2}\}$$

$$\text{Min}(2,2) = 2$$

$$\text{Min}(2,5) = 2$$

③ Sample Max

$$\text{Max} = \max \{\text{Num1}, \text{Num2}\}$$

$$\text{Max}(2,2) = 2$$

$$\text{Max}(2,5) = 5$$

④ Sample Range

$$\text{Range} = \text{Max}(\text{Num1}, \text{Num2}) -$$

$$\text{Min}(\text{Num1}, \text{Num2})$$

$P(\text{sample mean} \leq s) ?$



$P(\bar{X} \leq s)$