

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) Does 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 $i = 1, 2, 3, \dots, n$.
- (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 → $\{x_1, x_2, \dots, x_n\}$
 ↓
 Random sample
 ↓
 a collection of random variables
 ↓
 Changes every time 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
 ↓
 Each run of sampling procedure gets one of the elements in datasets

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

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

Sample data → $\{(x_1, x_2, \dots, x_n) \mid x_i \in \mathbb{R}, x_i = 1, 2, \dots, 5\}$

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

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

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

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