

Lecture - 01

08/26/2020

Survey: Who has a iPhone?

Standard notation for a "datum"

$X_1 = 0$; $X_2 = 0$, $X_3 = 1$, $X_4 = 1$, $X_5 = 0$, $X_6 = 1$, $X_7 = 1$, $X_8 = 0$,
first string "No" 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0
 X_{20}

$n = 20$ is our "sample" 12 1's ; 8 0's

Do we believe this survey is a "sample" of $n = 20$ elements from a superset called the "population"? If we do, this is called the "population model sampling assumption."

If so, what is that population?

- All people on Earth?
- All people in America?
- All college students?
- All college students in NYC?
- All public college students in NYC?
- All QC students?

Is this sample representative of the population?

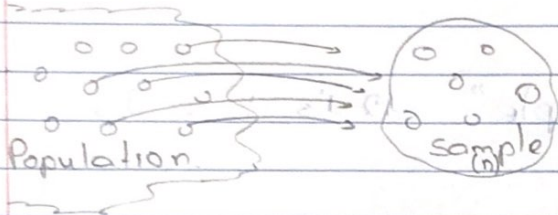
This is typical a sample, assume a population model, then identify the representative population. This happens in data science all the time.

In classical statistics, this goes in the opposite direction. You begin by defining the population

clearly and then sample n elements from that population.

Population has size N . You have some idea of what N is.

If population = all americans $\Rightarrow N = 330$ million



We see the data x_1, x_2, \dots, x_n in the sample but not other data in the population.

Can we learn about the population from the sample?

This is called "inference". We use the sample to "infer" properties about the population. Usually the properties are parameters of the random variable which creates the population.

"Infer" means to make an educated guess from specific things to universal properties. A synonym is "induction".

The opposite is deduction which is universal \Rightarrow particular. You can "NEVER" be sure your inference is correct.

How is inference done with data?

You generate "statistics" which are functions of the data.

$$\hat{\theta} = w(\overset{\text{data}}{x_1, \dots, x_n})$$

II \rightarrow statistic (usually scalar)

e.g. $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i = 0.6$
 $\hat{\theta} \equiv \bar{x} \text{ or } \hat{p}$ (our iphone survey)

What can you infer with this statistic?

usually, you infer θ , the population parameter which is the "true proportion" of iphones.

"Statistical inference" - using statistics to make inferences.

What is θ ?

$\theta = \frac{x}{N} \rightarrow$ # of people in the population that have iphones (unknown)
 $N \rightarrow$ # of elements in the population (known)
 parameter (unknown property)

$$\theta \in \Theta = \{0, 1/N, 2/N, \dots, (N-1)/N, 1\}$$
 the parameter space

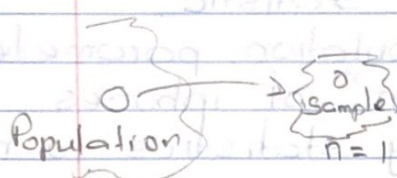
convention is that Greek letters represent unknown quantities and roman letters represent known quantities.

$\hat{\theta}$ is a "point estimate" for the unknown θ .
 "Point" meaning one specific value which you believe is a good guess for the value of θ .

I "Point estimation" is one of the goals of statistical inference. The other two are II Confidence set creation and III theory testing (testing a theory about a specific value of θ at a "certainly level" α .)

Let's sample one element from the population and do one survey.

How should this element be chosen if I want a "representative" sample?



Randomly but specifically, uniform meaning every element has probability of $1/N$ of being chosen. That's called a "Simple random sample" (SRS).

What is the probability that $X_i = 1$?

$$P(X_i = x_i = 1) = \frac{x}{N} = \theta$$

the r.v. \rightarrow modeling the Survey.
specific value \rightarrow the realization (a value in the support of x_i)