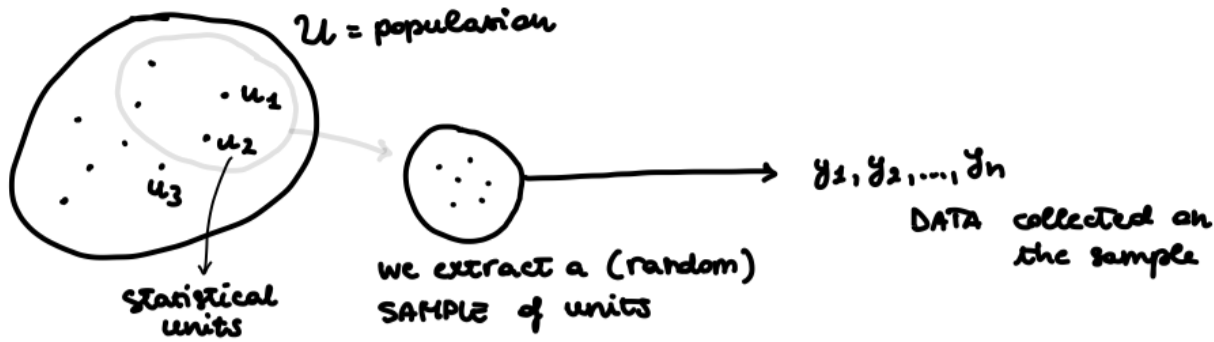


PREREQUISITES OF STATISTICAL INFERENCE



The observations are  $y_i = g(u_i)$

STATISTICAL INFERENCE: we try to understand the population starting from a sample

We use probabilistic tools

STATISTICAL MODEL: on the data, we specify a probabilistic model suited to describe a particular phenomenon.

probabilistic model:  $Y \sim p_\theta(y)$  we draw  $(y_1, \dots, y_n)$

statistical model: given  $(y_1, \dots, y_n)$ , we define a set of "reasonable" distributions  $P(y)$  that could have generated it, and try to recover the particular  $p_\theta(y)$  within this set

$y_i \rightsquigarrow Y \sim p(y; \theta) \quad \theta \in \Theta \subseteq \mathbb{R}^p$

known (prob. model)      unknown (parameter) identifies the particular element

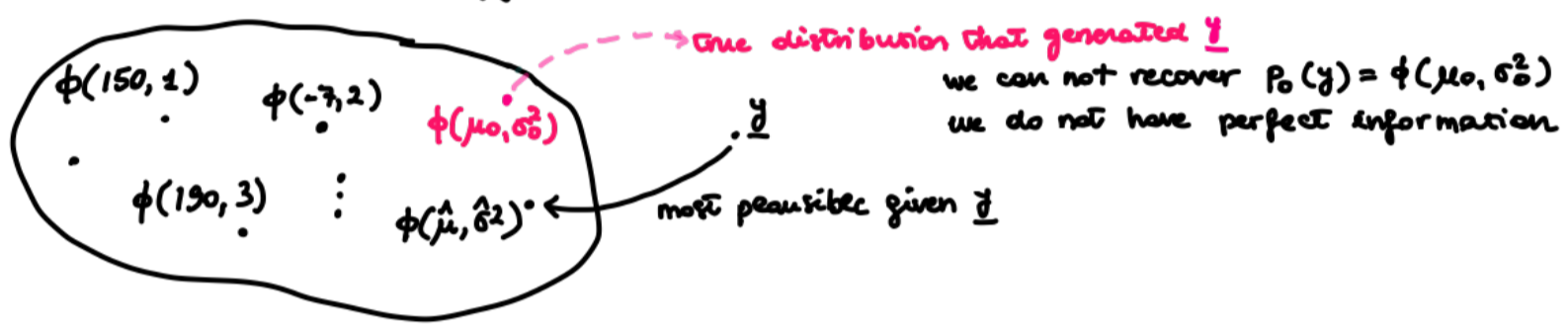
e.g.  $\underline{y} = (y_1, \dots, y_n)$  heights of a sample of students  
a Gaussian distribution is a reasonable model

$\Rightarrow$  statistical model  $Y_i \sim N(\mu, \sigma^2) \quad \mu \in \mathbb{R} \quad \sigma^2 \in (0, +\infty)$

random variable: before observing the data. The distribution of  $Y$  describes all possible realizations  $y_i$

MODEL: set of all Gaussian distributions

inference: we use the data to identify the element in this set that best describes them.



e.g.  $\underline{y} = (y_1, \dots, y_n)$  counts of cars passing in a street  
statistical model  $Y_i \sim \text{Pois}(\lambda) \quad \lambda \in (0, +\infty)$  indep.

Methodologies:

- POINT ESTIMATION: we identify one element (the most plausible) within the set of distributions  $\hat{p}(y)$
- INTERNAL ESTIMATION (confidence intervals): subset of reasonable elements, where the subset has a known degree of "confidence" that the true element ( $p_\theta$ ) is contained in it.
- HYPOTHESIS TESTING: we ask if there is enough evidence in the sample to draw conclusions about a particular statement ("null hypothesis")

POINT ESTIMATE

given the set of elements  $\{p(y; \theta) ; \theta \in \Theta \subseteq \mathbb{R}^p\}$  we want to identify the most plausible  $\hat{p}(y)$

It is identified through a particular element of  $\Theta$ :  $\hat{\theta}$  i.e.  $\hat{p}(y) = p(y; \hat{\theta})$

ESTIMATE  $\hat{\theta} = \hat{\theta}(y)$  is a function of the observed values (realizations)

ESTIMATOR  $\hat{\Theta} = \hat{\Theta}(Y)$  is a function of the random variable

$\Rightarrow$  we study the distribution of the estimator, its expected value, variance, ...

HYPOTHESIS TESTING

$\begin{cases} H_0: \theta \in \Theta_0 \subset \Theta & \text{null hyp.} \\ H_1: \theta \in \Theta \setminus \Theta_0 & \text{alternative hyp.} \end{cases}$  e.g.  $\begin{cases} H_0: \theta = \theta_0 \\ H_1: \theta \neq \theta_0 \end{cases}$

TEST: we partition the sample space into the REJECT ( $R$ ) and ACCEPTANCE ( $A$ ) regions

- $A$ : values of  $y$  that suggest that  $H_0$  is true
- $R$ : values of  $y$  that suggest that  $H_0$  is false  $\rightarrow$  reject  $H_0$

TEST STATISTIC: a function of the data that defines the two regions.  $T(y)$

- $A = \{y \in \mathcal{Y} : T(y) \text{ suggests } H_0\}$
- $R = \{y \in \mathcal{Y} : T(y) \text{ suggests } H_1\}$

How do we draw conclusions?

- FIXED SIGNIFICANCE LEVEL  $\alpha$   
we guard against the 1<sup>st</sup> type error: we fix  $\alpha$  to a small value (e.g.  $\alpha = 0.01, \alpha = 0.05, \alpha = 0.10$ ) and we derive  $A$  and  $R$  so that  $IP(\text{reject } H_0 \mid H_0 \text{ is true})$  is equal to  $\alpha$  (or at most  $\alpha$ )  
In other words,  $\alpha = IP(\text{1<sup>st</sup> type error}) = IP(y \in R \mid H_0 \text{ true})$   
Of course, the smaller  $\alpha$  is, the smaller  $R$  will be (I want to reject  $H_0$  only if I am really really confident)
- OBSERVED SIGNIFICANCE LEVEL (p-value)  $\alpha^{obs}$   
it is the probability of observing "more extreme" values (i.e. more against  $H_0$ ) than the ones we observed.
  - if the reject region is of a one-tailed test
    - $H_1: \theta > \theta_0$  (right tail)  $\Rightarrow \alpha^{obs} = IP_{\theta_0}(T \geq t^{obs}) = IP(T(Y) \geq t(y^{obs}) \mid H_0 \text{ true})$
    - $H_1: \theta < \theta_0$  (left tail)  $\Rightarrow \alpha^{obs} = IP_{\theta_0}(T \leq t^{obs})$
  - if the reject region is of a two-tailed test
    - $H_1: \theta \neq \theta_0 \Rightarrow \alpha^{obs} = 2 \min\{IP_{\theta_0}(T \geq t^{obs}); IP_{\theta_0}(T \leq t^{obs})\}$

The two procedures are related: if  $\alpha^{obs} < \alpha$ , then I reject  $H_0$  in a fixed-level test of level  $\alpha$

CONFIDENCE INTERVALS of confidence  $(1-\alpha)$

it is a random interval  $\hat{C}(Y)$  such that  $IP(\theta \in \hat{C}(Y)) = 1-\alpha$  for all  $\theta \in \Theta$

With probability  $(1-\alpha)$ , the interval contains the true value of the parameter, whatever it is.

After we compute the interval with the data (hence, we get a fixed numeric interval), it either contains the true  $\theta$  or not.

The probability must be interpreted regarding to the random quantity.

We build it through the identification of a PIVOTAL QUANTITY: a function  $g(Y; \theta)$  of the r.v.  $Y$  and the parameter  $\theta$  such that its distribution does not depend on  $\theta$  (hence it is completely known).

Then we look for the interval  $(u, v)$  such that  $1-\alpha = IP(u < g(Y; \theta) < v)$ .

With the data we compute  $\hat{C}(y^{obs}) = \{\theta \in \Theta : g(y^{obs}, \theta) \in (u, v)\}$

$\hookrightarrow$  all values of the parameter that, given the observed data, give a value of  $g$  within the  $(u, v)$  interval.