

Lecture 4. Intro to Statistics

Basic Concepts: Population, Sample, Parameter, Statistics

so make mistakes all the time

is one realization from pop. // excel table

- Sample= a single draw of data from all potential realizations of that data (random vector \mathcal{X})
- Population= the distribution F_X of data vector \mathcal{X} .
- Probability Theory: given population what can you say about distribution of $g(\mathcal{X})$?
 \rightarrow physics
- Statistics: ~ engineering
given a single realization of \mathcal{X} what can you say about F_X ?

Basic Concepts: Population, Sample, Parameter, Statistics

3 types of data:

- Cross-section* = a set of i.i.d. random vectors
 $(\mathcal{X} = (X_1, \dots, X_n), x = (x_1, \dots, x_n))$. If we assume that $X_i \sim F$, then
 $F_X(x) = \prod_{i=1}^n F(x_i)$.
- Time-series* $X_t, t = 1, \dots, T$ allows dependence between consecutive observations . $\mathcal{X} = (X_1, \dots, X_T)$ is one realization of a path of a dynamic process.
- Panel data* $\mathcal{X} = \{X_{it}, i = 1, \dots, n, t = 1, \dots, T\}$ a draw from n independent identically distributed dynamic processes.

Basic Concepts: Population, Sample, Parameter, Statistics

→ fixed, known
not random
unknown & trying to learn

- **Parameter** = functional of the unknown distribution F_X (population concept!). $p = \text{prob of success}$
- **Inference** = to render a judgement about a parameter (or population F_X) based on a single draw (sample).
- 3 types of inference: estimation, confidence sets, testing

Ex. 1) coin toss $x_i = \begin{cases} 1 & p \\ 0 & 1-p \end{cases}$
 $\underbrace{x_1, \dots, x_n}_{\text{sample}} : \text{iid}$ $p: \text{parameter}$
 $\underbrace{\text{sample } [x_1, \dots, x_n]}_{\text{based on one realization infer } p}$

Ex. 2) $(x_i, y_i) \sim \text{iid } F_{xy}$
 $x_i \leftarrow \text{height}$
 $y_i \leftarrow \text{income}$

one draw $\leftarrow \{(x_1, y_1), \dots, (x_n, y_n)\}$
excl. sample

i think θ exists

$\theta = \text{cov}(x_i, y_i)$
so look at
sample and
say something
about pop

Basic Concepts: Population, Sample, Parameter, Statistics

- Statistic = function of sample $Y = g(\mathcal{X}) = g(X_1, \dots, X_n)$
- It is a random variable! *so interested in its dist'n & what we can tell about it*
- The distribution of a statistic = sampling distribution
how it appears from different samples

Sample Mean and Sample Variance

"average"

□ **sample mean** ($\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$)

□ **sample variance** ($s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$)

$$\left. \begin{array}{l} E[\bar{X}_n] = \mu \\ \text{mean of } \bar{X}_n \text{ is } \mu \\ \text{Var of } \bar{X}_n = \sigma^2/n \\ \text{Var}[\bar{X}_n] \\ \text{Var}(X_i) = \sigma^2 = \\ = E(X_i - \mu)^2 \end{array} \right\}$$

Lemma 1.

If $X_1, \dots, X_n \sim i.i.d.F$ is a random sample of size n from a population distribution with mean $\mu = \mathbb{E}X_i$ and variance $\sigma^2 = \text{Var}(X_i)$, then

$$\mathbb{E}[\bar{X}_n] = \mu \text{ and } \mathbb{E}[s^2] = \sigma^2.$$

$$\left. \begin{array}{l} \mathbb{E}\bar{X}_n = \mu \\ \text{pop; parameter; unknown} \\ \text{can be calc from sample} \\ - - - \\ s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad \& \text{wts mean of } s^2 \text{ i.e. } E s^2 = \sigma_n^2 \\ (\text{sample var}) \quad \text{h.v.} \end{array} \right\} \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

Empty slide

$$\square \quad y_i = x_i - u \quad \text{var}(y_i) = \sigma^2 \quad \bar{y}_n = \bar{x}_n - u \quad y_i - \bar{y}_n = x_i - \bar{x}_n$$

$$\text{so} \quad s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \quad : \quad \frac{1}{n-1} \sum_{i=1}^n [y_i^2 - 2\bar{y}y_i + \bar{y}^2]$$
$$= \frac{1}{n-1} \sum_{i=1}^n y_i^2 - 2\bar{y} \sum_{i=1}^n y_i + n\bar{y}^2$$

$\underbrace{n\bar{y}}_{n\bar{y}}$

$$s_n^2 = \frac{1}{n-1} [\sum y_i^2 - n(\bar{y})^2]$$

$$E y_i = 0$$

Taking Exp,

$$E s_n^2 = \frac{1}{n-1} [\sum_{i=1}^n E y_i^2 - n E (\bar{y})^2]$$

\downarrow

$\text{var}(y_i) = \sigma^2$

$\text{var}(\bar{y}) = \frac{\sigma^2}{n}$

$$= \frac{1}{n-1} [n \cdot 6^2 - n E(\bar{x})^2] = \sigma^2$$

$E(x-a)^2 \rightarrow \min_a$

$\frac{1}{n} \sum (x_i - a)^2 \rightarrow \min_a$

$a = \frac{1}{n} \sum x_i = \bar{x}$

But \bar{x} is pop. if given from outside

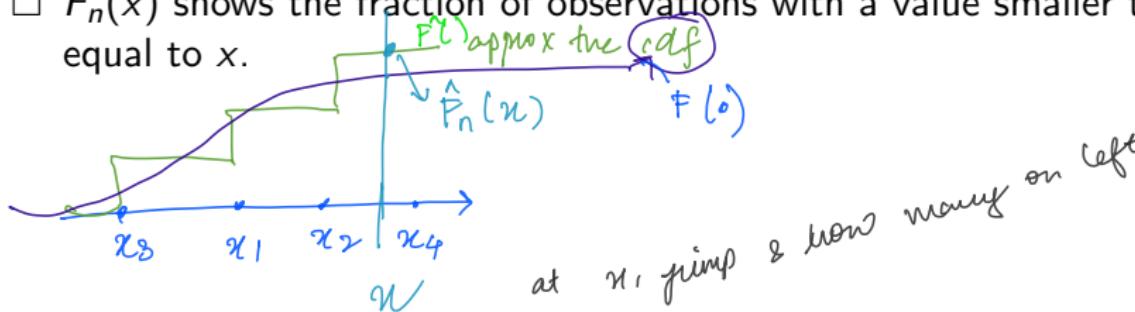
$E \frac{1}{n} \sum (x_i - \bar{x})^2 = \sigma^2 // \text{random}$

$\frac{1}{n} \sum (x_i - u)^2 \geq \frac{1}{n} \sum (x_i - \bar{x})^2$

Empirical distribution function

\hat{F}_n approx est.

- \hat{F}_n is the cdf of the distribution that places mass $1/n$ at each data point X_i .
- $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(X_i \leq x)$,
- $\hat{F}_n(x)$ shows the fraction of observations with a value smaller than or equal to x .



Empirical distribution function

LLN: Law of large no's

Theorem 1.

If we have a random sample X_1, \dots, X_n of size n from a distribution with cdf F , then for any $x \in \mathbb{R}$, $\mathbb{E}[\hat{F}_n(x)] = F(x)$ and $\text{Var}(\hat{F}_n(x)) \rightarrow 0$ as $n \rightarrow \infty$. As a consequence, $\hat{F}_n(x) \rightarrow_p F(x)$ as $n \rightarrow \infty$.

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i \leq x\}}$$

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n y_i = \bar{Y}$$

$$\mathbb{E} \hat{F}_n(x) = \mathbb{E} \bar{Y} = \mathbb{E} Y_i = p = F_X(x)$$

$$\text{LLN: } \hat{F}_n(x) \xrightarrow{P} F_X(x) \quad \text{as } n \rightarrow \infty$$

| my fn at every point does converge to cdf

a fn converges to fn

(both monotonic not
have much)

Empirical distribution function

Space to
depart from
each other)



Theorem 2 (Glivenko-Cantelli).

If X_1, \dots, X_n is a random sample from a distribution with cdf F , then

$$\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \rightarrow_p 0.$$

Ways to find the distribution of a statistic

analytical,
rarely used.

Method 1: exact distribution

- Question: If you know the distribution of \bar{X} , what is the distribution of $Y = g(\bar{X})$?
- Can you do this analytically?
- Very rare case.

Ex 1 $x_1, \dots, x_n \sim \text{iid } N(\mu, \sigma^2)$
 $Y = \bar{x} \sim N(\mu, \sigma^2/n)$ -

Ex 2 $x_1, \dots, x_n \sim \text{iid}$ $x_i = \begin{cases} 1 & P \\ 0 & 1-P \end{cases}$ || Bernoulli
 $Y = \sum_{i=1}^n x_i$ $P\{Y=k\} = \frac{n!}{k!(n-k)!} P^k (1-P)^{n-k}$ || Binomial
(y can take values $1, 2, \dots, n$)

Ways to find the distribution of a statistic :

(simulation)

Method 2: Monte-Carlo (Numerical way to implement Method 1)

- Question: If you know the distribution of \mathcal{X} , what is the distribution of $Y = g(\mathcal{X})$?
 $x_1, \dots, x_n \sim \text{iid } F$
- For $b = 1, \dots, B$, simulate $\mathcal{X}_b^* = (X_{1b}^*, \dots, X_{nb}^*)$, where $X_{ib}^* \sim \text{i.i.d } F$
- Calculate $Y_b^* = g(\mathcal{X}_b^*)$. calc. value of statistic; give multiple draws y_1^*, y_2^*
- We can approximate distribution of Y :
 - $F_Y(s) \approx \frac{1}{B} \sum_{b=1}^B \mathbb{I}\{Y_b^* \leq s\}; \rightarrow y_1^*, \dots, y_B^*$
 - $P\{Y \in A\} \approx \frac{1}{B} \sum_{b=1}^B \mathbb{I}\{Y_b^* \in A\};$
 - α -quantile of Y : $q_Y(\alpha) \approx Y_{(\lfloor \alpha B \rfloor)}^*$; $y_{(1)}^* \leq y_{(2)}^* \leq \dots \leq y_{(B)}^*$
mean of y by taking avg of y
 - $\mathbb{E}Y \approx \frac{1}{B} \sum_{b=1}^B Y_b^*$; ordered from L to R
rounding up or down // took α of them.
 - $\text{Var}(Y) \approx \frac{1}{B-1} \sum_{b=1}^B \left(Y_b^* - \frac{1}{B} \sum_{s=1}^B Y_s^*\right)^2$.

You control the quality of approximation!

Q) How do we simulate from same distn as pop
v.good simulation; if control B can do as many & get good-

Empty slide



Ways to find the distribution of a statistic

CLT: start w/ wholen
distribution — stabilize
& get normal

Method 3: Asymptotic approximation (most often used method)

- What if you do not know the distribution of \mathcal{X} ? Is it helpless to figure out the distribution of $Y = g(\mathcal{X})$? $x_1, \dots, x_n \sim \text{iid}$
- You can approximate distribution of some statistics by the limit as $n \rightarrow \infty$
- Looking for statement

$$F_{Y,n}(t) \{ Y = g(\mathcal{X}_n) \leq t \} \Rightarrow F_\infty(t). ?$$

- We rely on CLT, delta-methods, Slutsky theorem - to get a distribution that can be normal.
- We do not control the quality of approximation (depends on n and speed of convergence) ↓

Empty slide



Ways to find the distribution of a statistic

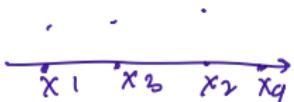
simulate data from my empirical cdf which is actually close to the true/real cdf of pop.

Method 4: Bootstrap*

- What if you know the distribution of \mathcal{X} ? Is it helpless to figure out the distribution of $Y = g(\mathcal{X})$?
- Idea: to approximate F_X by some close distribution.
- If Y depends *continuously* on F_X ,
and if $F_X \approx \hat{F}$, then we approximate Y
- The last step can be done by Monte-Carlo.

empirical cdf:
draw from my sample.

I have a sample



Ways to find the distribution of a statistic

① my data comes from empirical distⁿ;

non-parametric

Method 4: Bootstrap Glivenko-Cantelli' theorem: $\hat{F}_n(x) \approx F(x)$. can control

- For $b = 1, \dots, B$, simulate $\mathcal{X}_b^* = (X_{1b}^*, \dots, X_{nb}^*)$, where $X_{ib}^* \sim \text{i.i.d } \hat{F}_n$ simulation draw X_{ib}^* independently with replacement from the set $\{x_i, i = 1, \dots, n\}$.
- Calculate $Y_b^* = g(\mathcal{X}_b^*)$.
- We can approximate distribution of Y by that of Y_b^* :

- $F_Y(s) \approx \frac{1}{B} \sum_{b=1}^B \mathbb{I}\{Y_b^* \leq s\}$;
- $P\{Y \in A\} \approx \frac{1}{B} \sum_{b=1}^B \mathbb{I}\{Y_b^* \in A\}$;
- α -quantile of Y : $q_Y(\alpha) \approx Y_{(\lfloor \alpha B \rfloor)}^*$;
- $\mathbb{E}Y \approx \frac{1}{B} \sum_{b=1}^B Y_b^*$;
- $\text{Var}(Y) \approx \frac{1}{B-1} \sum_{b=1}^B \left(Y_b^* - \frac{1}{B} \sum_{s=1}^B Y_s^* \right)^2$.

You do not control the quality of approximation.

② more here is emp. cdf to true comes from org. n. long sample $\mathcal{D}(n)$ cannot control

Empty slide



Ways to find the distribution of a statistic

Plug-in approach

- A random sample X_1, \dots, X_n from a population F
- Want to estimate $\theta = T(F)$ (functional of population)
- Plug-in idea: if $\hat{F}_n(x) \approx F(x)$, then $T(\hat{F}_n) \approx \theta$.
- This is informal idea.
- Example

$$\mu = \mathbb{E}X_i = \int x dF(x) \approx \int x d\hat{F}_n = \bar{X}_n$$

if you know
cdf
mean
mean

$$\int x d\hat{F}_n(x)$$

puts equal weight on all obs

$$= \frac{1}{n} \sum x_i = \bar{x}$$

$$m = \text{med } x_i = F(m) = 1/2$$

try to find \hat{m} s.t.

$$\hat{m} : \hat{F}_n(\hat{m}) = 1/2$$
$$X([m/n])$$

Empty slide

sample

$$\text{Var} = \int (x - \text{SredP})^2 dF$$

□ plug in Var $\hat{\text{Var}} = \frac{1}{n} \sum (x_i - \bar{x})^2$

j