

Statistics

- Statistics = The science of learning from data
- It gives us methods & tools to:
 - collect data
 - summarize data
 - Analyze data
 - interpret result
 - make decisions under uncertainty

Statistic vs Statistics

→ Statistic → A single number or summary from data

Ex → avg exam score, minimum Temp, % of retired people in a survey.

• Statistics → Academic discipline

→ Develops methods to analyze, interpret, and apply statistics

→ works with experts to make meaningful conclusions.

why Statistics matters:

(1) summarizing data: → makes huge data sets understandable.

Ex → average, graphs, percentiles.

(2) Avoiding misleading info → Provides structures to check if claims are valid.

Ex → Opinion polls with margin error

(3) Decision making under uncertainty

Ex → Should a patient undergo preventive treatment if at risk?

(4) Understanding variation → not just average, but spread of data.

Ex → weather, election, demand for a product

(5) Prediction / forecasting

→ weather, elections, demand for a product

6. measuring difficult concepts

- Easy: age, height
- Hard: mood, Political ideology.
- Statistics helps design rigorous measurement methods

7. Efficient data collection

- Balance btw too much data (wastly) vs too little data (inaccurate)

where Data come from

→ why does Data source matter?

→ knowing how data is generated is crucial
it influence which statistical methods are valid.

→ Before analyzing → ask

- where did the data come from
- was it designed or organic?
- can it be assumed i.i.d?

Two main types of data

I) organic / Process data

- Data generated naturally by ongoing processes.
- usually large-scale (Big data) & collected automatically
- Example:
 - financial transaction, stock markets
 - netflix viewing history
 - web browsing activity
 - sports performance metrics
 - sensor data (Temp, pollution)

Features:

- ~~probabilistic~~ → massive in size → requires computational resources
- Often used to find patterns & trends
- Harder to clean & prepare for analysis

(2) Designed data collection

- Data gathered from Purposful Studies to answer specific questions.
- Example
 - surveys (opinions from a sample of people)
 - interviews with sampled individuals
 - extracting and coding specific tweets for research.

features

- smaller in size
- easier to work with.
- designed with a clear research objective
- more controlled & rigorous than organic data

Big data

- refers mainly to organic process data
- challenges
 - Storing, cleaning, processing.
 - Required stronger + statistical methods to make usable datasets

i.i.d. data (independent & identically distributed)

- independent → each observation doesn't affect others
 - identically distributed → all observations come from the same statistical distribution
- Ex → final exam score in a large stats class - each student grades independently, score may follow a normal distribution

→ why i.i.d. is important

→ many statistical method assume i.i.d.

→ Allows estimation of:

→ mean, variance, Percentile

→ Precision of these estimates

→ Non-i.i.d. cases (when assumptions break)

→ dependence: students cheat → scores not independent

→ Difference distributions: males vs females may have different score distributions.

→ group effects: - students in different discussion sections may perform differently

in such cases

→ must use different statistical methods to handle dependencies.

→ Always ask:

1. did data come from organic process or designed collection?

2. can we assume i.i.d?

→ choice of statistical procedure depends on these answers.

Variable type

→ why variable matters?

→ The type of variable determines:

→ how we summarize it (mean, frequency)

→ what analysis methods make sense

NHANES

NHANES DATA

ID	BMI	Race	Age	Adult
62161	23.3	3	22	1
62163	17.3	5	14	0
62164	23.2	3	44	1
62165	27.2	9	16	0
62202	24.7	1	36	1
...

Adult

0 - Age is less than 18

1 - Age is greater than equal to 18.

Two main types of variable

(1.) Quantitative (Numerical):

→ Represent measurable quantities

→ Arithmetic (like average, sum) makes sense

Subtypes:

→ continuous → can take any value within a range

ex → BMI, height, weight, time to run a mile, age

→ discrete → only specific countable numbers

→ ex → no of children in a household
numbers of pets

Race

1. Mexican America
2. Other Hispanic
3. Non-Hispanic white
4. Non-Hispanic Black
5. Other

2. Categorical (Qualitative)

- classifies items into groups
- Arithmetic doesn't make sense
(and average race codes meaningless)

Subtypes:

- Ordinal → has a natural order/ranking
 - Ex → class standing (freshman → sophomore → junior → senior)
- Nominal → no natural order
 - Ex → race, marital status, eye color etc.