# Getting Started with Statistics (Day-06)

↳ Statistics is the Science of collecting, **Organizing** and
— analyzing data.
&
Decision Making

Data ——| "facts or pieces of info"
( measured, Calculated, Analyzed )

Eg:-
① weights of students in a class

$$\{ 60, 50, 45, 30 \cdots \}$$

② IQ of the students in a class

$$\{ 100, 90, 95, 80 \cdots \}$$

## Data set :- House price Dataset

| City | Area | No. of rooms | Price | |
|------|------|--------------|-------|---|
| Bangalore | 1000 | 2 | 45Lakhs | Analyze Date |
| New york | 1250 | 25 | 50 lakhs | Data Scienti — 1MO |
| mumbai | - | - | - | ↓ Price |

Data Analyst ——| Report ——| Visualization ——| meaning Decision
↳ project —|

# Application

1. Data exploration & Summarization
2. Model building & validation
3. Statistical Analysis ——| Sample data —| population data.
4. ~~Hypothesis~~ testing
5. optimization & Efficiency.
6. Reporting.

# Types of Statistics

1. Descriptive Statistics

   └ It involves methods for Summarizing and Organizing data to make it understandable.

   └ This type of Statistics helps to describe the basic features of the data in a study.

   ① measure of To Central tendency
   (mean, median, mode)

   ② measure of Dispersion ( variance, standard deviation.

   ③ Data distribution
      (i) Histograms
      (ii) Box plot
      (iii) Pie chart
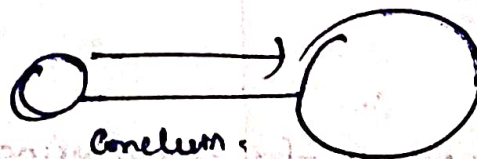      (iii) PDF, PMF

④ Summary statistics

① Five Number summary

$Q_1, Q_2, Q_3$, Max value.

② Inferential Statistics -

└ It involves methods for making predictions or inferences about a population based on a sample of data. It allows for hypothesis testing, estimation & drawing Conclusions.

Sample data ⟶ ◯ ━━➤ ⬭
                    Conclusn.

① Hypothesis testing

② P value

③ Confidence Interval

④ statistical Analysis test
   ① z-test
   ② t-test
   ③ F- test (ANOVA)
   ④ chi square test

Eg:- Let Say there are 20 statistics class in your college
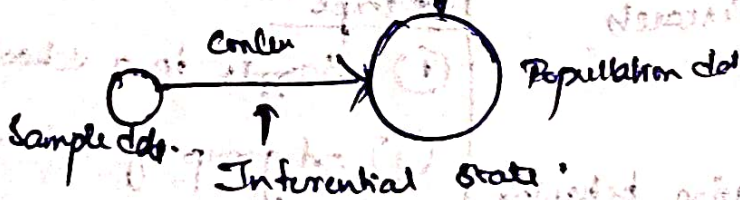and you have collected the height of students in the
class.

Heights are recorded { 175cm, 180cm, 140cm, 135cm, 165cm, 120cm

## Descriptive statistics

" what is the Avg height of the entire class "
or
measure of central tendency.

## Inferential statistics !

" Are the height of sample student in class similar to what you
expect in the entire college " ?



Sample data ─ Conclu → Population dat
                Inferential stats

* population and sample data
     ↓
* It is the entire set of objects of
  interest in a particular study.

* A Sample is a subset of the
  population that is used to represent
  the entire group.

| Population | Sample data |
|---|---|
| **Characteristics :** | **Characteristic** |
| 1) Complete set | ① Subset of population |
| 2) Parameter | ② Statistic |
| → A Numerical value summarise the entire the populat [Eg : population mean, population variance] $(\mu)$ $(\sigma^2)$ | → A numerical value summarise by Sample data ① Sample mean ② Sample variance |
| 1) population in a school study ● All students enrolled in a school) ↓ Avg height of student ↑ Population mean | ③ Random Sampling ! Sample should be randomly selected , to avoid bias, |
| 2) population in market Research All Consumer in a city To understand the purchasing behaviour of all Consumer, | **example** ① Sample in a school study ① A group of 50 student from school ② estimate the Avg height of student in a school. |
| 3) population in a medical study ● All the patient with a specific desear ● To study the effectiveness of a drug, | ② Sample in market rescarch ② 500 Consumer from the city relai - Behaviour ——→ population ③ sample medical study 200 patient use ca Test the effectiveness of the drug |

Types of sample Techniques :-

① Po' probability sampling

② simple Random sampling

⤷
→ Every member of the population will has an equal chance of being selected.

Eg:- Selecting people randomly.

Draw names random from a class of students.

ⓑ ~~Systma~~ Systematic sampling :

⤷ select every $n^{th}$ member of the population after a random starting point.

Eg:- Airport ——+ Credit card ——— $5^{th}$ person, $10^{th}$ person, $15^{th}$ person.

Feedback survey !

ⓒ Stratified sampling :

Divide the population into Strata (groups) based on the specific characteristic & then randomly sampling from each Group.

Eg:- Divide employees by department in a Company & select a proportional number from dept to from a survey Sample.

Eg → Age ————— < 12     12 - 18    > 18 d country

① Cluster sampling

   ↳ Divide the population into clusters, randomly selecting

   Clusters, then sampling all the members from the selected

   Clusters.

Eg → randomly selecting several schools from a district, and

     Surveying all trade within those schools

② Multi stage Sampling -

   Combining several sampling method, usually involve

   Selecting cluster, then randomly sampling within the

   Cluster

Eg → randomly selecting City, each selected city randomly

   selecting household to survey.

① **Non - probability Sampling**

Ꮮ₁ select individual who are easiest to reach

Eg:- Surveying people in the mall

⑳ **Convenience Sampling:**

Ꮮ₁ selecting individual who are easiest to reach

⑱ **Judgemental Sample**

Ꮮ₁ select individual ~~who are~~ based on the researcher's Judgement

↑useful for representation.

Eg: choose expert in a field to participate & Datascience y

⑳ **Snowball Sampling:**

existing study subjects future subjects from among other. requirements,

⑤ Types of Data

① Quantitative        ② Qualitative

| Discrete | Continous | Nominal | Ordinal |
|---|---|---|---|
| ↓<br>① whole Number<br>Egt No. of bank<br>  accounts<br>No. of children in a<br>family | * Any value<br>Egt weight, height,<br>temputure, speed, | Eg: gender (M,F)<br>bloodgroup.<br>Pincode<br>→ Categorical value,<br>→ No rank | Eg: Consumer feedbk<br>Good, bad, best<br>→ Contains rank |

Scales of measurement of Data
Lı Describes the nature of into within the value assigned to
   ~~measurement~~ Variable,

① Nominal Scale
   Lı This scale classifies the Data into Distinct Categories
   that do not have an intrinsic order,
   Qualitative / Categorical data,

characteristics
① Categorized based on labels, names r qualitu.
② Categorized are mutually exclusive.

(ii) No logical order among category [No rank]

Ex:- Gender       Colour
    → M       red — 5      50%
    → F       Blue — 4,    40%
              Pink — 1     4%6
                    ——
                    ∞

★ Ex:- Types of cuisines
    { Italian }
    { chinee  }
    { media   }

② **Ordinal Scale :**

   └ Classifies the data that can be ranked or ordered,

Characters

   ① Have logical order among category [rank]

   ② Interval b/w rank are not necessarily equal.

Ex:                                    rank
    Educational level      =    If    —1
       High School                     1
       Bachelor                        2
       masters                         3
       Doctorate                       4

Or: Customer feed back ,    Very satisfied 2  , unsatisfied  0
    Satisfied —1      ,

③ Interval Scale :

└ The interval scale not only categories and order but also specify the exact distance, blw interwale.

→ st. lost lacks a true zero point.

Characteristics:

① ordered with consistent interval blw value.

② Allows for meaningful comparision of difference.

③ No true zero point.

Enter

Temp in Farenheit

$10°F$ , $20°F$ , $30°F$

$=1 8F$ ~1 zero temp

ST $\phi$ sem :

$90, 100, 110$

LDST $100 - 90 = 10$

$\therefore$ STp $\neq 10$

Calender year

$2024, 2020, 2016$ ≠1 0 yr.

(4) **Ratio Scale** :

* Order matters

* Differences are measurable ( ratio can be measured,

* Contains a 0 starting point

Eg: marks of the students in a class.