

SDS → Unit 1

- Data science → Interdisciplinary field that extracts insights from data (ie applies to solve problems across various domains)
- Applications → Image recog., speech recog, internet search, digital ads, price comparison, fraud & risk detection.
- Data → Individual facts / Items of info collected through an obsy.

Data → raw facts formatted in special way. → Based on records, observations & unorganized.

Information → Data org. in such a way that they have addⁿ, value beyond the facts themselves.

DATA, ENTERPRISE SYSTEMS

- Structured data → Data whose elem are addressable for eff analysis. stored in formal repos, → database. Ex: Relational data.
- Semi-structured data → Doesn't reside in database but has some org. prop. that make it easier to analyse. Ex: XML Data.
- Unstructured data → Not org. in predefined manner or have data model! Ex: word, pdf
- Knowledge → Body of guidelines & procedures to select/org./manipulate data to perform specific tasks on data.
- Science → Enterprise that org. knowledge to form testable expr. & pred. of the universe.

- Six Vs of BigData →
- Volume → amt of data collected
- Variety → types of data generated
- Velocity → speed at which data is generated
- Veracity → trustworthiness of data
- Value → business value of data
- Variability → ways of using/formatting data

→ Using data →

→ Explore → Qd, patterns

→ Predict → make informed guesses

→ Inference → quantify what you know
bring data to model, further from prod, data given

→ Statistics → use math, eqn, & stats model to
analyse data & make conclusions.

→ Data scientist → deliver actionable results
/ deploy model to the prod, system.

bring
model to data,
embedded in prod,
Need to get data

→ Data analysis → collecting, exploring large amt of
data to discover underlying patterns & trends.

→ Pop. → entire collection of obj / outcomes
from which info is sought.

→ Sample → subset of pop. → actually obj
outcomes / obj.

- Sampling → process of selecting obj. in order to make inference generalised to pop.
- Reasons for sampling →
 - necessity → can't study whole pop.
 - practicability → more eff.
 - cost eff.
 - manageability → small datasets easier
 - time saving
- Char of sample →
 - rep. of pop.
 - app. size
 - unbiased sample
 - randomly selected
 - economical
 - goal oriented, std. systems
- Tangible pop. → members are phys. obj., always finite & involves counting.
- Conceptual pop. → not phys. obj., result of measure, not well defined, size → large
- Target/theoretical pop. → entire grp. on which researcher interested in generalising conclusions.
- Study accessible pop. → grp. on which researcher actually apply conclusions to.
- Sampling frame → list of items/events from which pop. respondents drawn to measure. Must be rep. of pop.

- Prob. Sampling → every unit in pop has chance/prob. of being selected → Equal prob. selection (EPS) / self-weighting.
- Non-prob. Sampling → not EPS. Selection based on assumptions on pop. of interest. Likely to produce biased sample & restricts generalisation
- Simple random sampling →
 - EPS, selection using random no. system/lottery.
 - Sampling frame → whole pop.
 - Random → repr. of pop., used when pop. ↓
 - Small pop. → lottery, Large → system gen. no.
- Prob. = $\frac{n \times 100}{N}$
 - $n \rightarrow$ Sample size
 - $N \rightarrow$ Pop. size
- Adv., →
 - Easy to use calculate
 - EPS, repr. of pop.
 - ↓ sampling error
 - min. knowledge of pop. needed in adv.
- Disadv., →
 - not useful when pop. ↑
 - may not repr. min. subgrp.
 - Can't be used when units → heterogeneous

→ Systematic Sampling →
→ arr, pop, accd to ordering scheme & select elem at reg. intervals. First elem → round off.

→ Select every k^{th} element, $k = N/n$.

→ EPS. Not simple random → diff. subsets, diff. prob.

→ Adv →

→ easy to select, sampling frame by easily

→ cost eff., entire pop, evenly sampled

→ ↓ chance of data contamination.

→ Disadv →

→ might lead to bias.

→ diff. to assess precision of estimate.

internally homogeneous → ignore of elem b/w k^{th} elem

→ Stratified Sampling → based on same shared char

→ Pop. div. into 2 grp, → strata & then sample using simple random.

→ Don't overlap → rep. whole pop.

→ Adeq, rep. of many subgrp.

→ Adv →

→ ↑ rep., easy

→ ↑ stat. off, ↑ precision

→ Disadv →

→ time consuming, ↑ cost

internally heterogeneous → classification errors

→ One stage

→ Cluster sampling → Two stage

→ pop. div. into non-overlapping clusters & then each cluster → randomly selected.

→ Heterogeneous members on each grp.

→ Adv →

→ useful for geo, diverse pop,

→ ↓ travel cost, ↓ resources, ~~more feasible~~

→ more feasible, administration easier.

→ Disadv →

→ ↓ eff, ↑ cost than simple randoms

→ ↑ sampling error, prone to bias.

→ Convenience sampling →

→ Sampling from pop, close to hand/convenient to researcher.

→ Adv → ↓ cost, ↓ time, simple, easy to imp,

→ Disadv → bias, not rep. of pop, sampling error

→ Judgemental Sampling →

→ Researcher chooses sample based on who they think is app for the study. Used when limited no. of ppl that have expertise in area being sampled

→ Adv → ↓ time, real time results, no special knowledge of stats needed

→ Disadv → prone to error, ↓ reliability ↑ bias, can't generalise findings

→ Quota Sampling →

→ Split into mutually exclusive grp, & then judgement applied on each based on predetermined characteristics of population.

→ Adv → ↓ cost, fast, convenient, rep. of any grp

→ Disadv → impossible to det. sampling error, bias, not possible to make stat. inf.

- Snowball sampling →
 - Survey subj, from referral from other survey subjects
 - Existing subj, asked to nominate further subjects
 - Used when sampling frame diff to pop.
 - Adv → ↓ cost, simple, little planning, ↓ workforce
 - Disadv → bias, impossible to def^s sampling error
rep, not guaranteed

- Sample statistic → piece of info, from fraction of pop.
- Pop ~~statistic~~ parameter → "refers to whole pop."
- Sampling errors/random errors →
 - discrepancy b/w sample stat & pop param
 - Ocurs when sample not rep. of pop
- Non-sampling errors/systematic error →
 - Result of errors made during data collection & processing → mostly due to sampling bias & non-response bias
 - ↓
 - when sample not rep. of pop

- Selection bias → Intended grp, ↑, ↓ prob, than others
- Non-response bias → sampling bias due to absence of certain obj/grp
- Sampling variation → when two diff samples from same pop, differ from each other. Hence
- Independence → knowing value of some cant helps predict values of other items

- Data → Qualitative → Nominal → Ordinal
- Data → Quantitative → Interval → Ratio
- Quantitative → recorded on numeric scale.
- Discrete Continuous
- Qualitative → not recorded on numeric scale but in categories
 - Nominal Ordinal
 - Labels for categories w/o direction/order. Ex: Male, female
 - Can't be manipulated with math ops.
 - Setting up inequalities, but no abs value
 - Labels for categories with order. Ex: first, second, third
- Ordinal data → order matters but not difference b/w them
 - Ex: $5 < 7$ but diff b/w 3 & 5 ≠ diff b/w 5 & 7.
- Discrete → limited no. of possible values & jumps b/w data. Ex: Integer nos.
- Bar graphs used to rep. discrete data.
- Continuous → can take any value within an interval measured on scale or continuum.
- Interval → data measured on scale with each pt. equidistant from each other. Measurable, ordered but no meaningful zero. Ex: Mean, Median, Mode, min, max, temp.
- Ratio → classifies & ranks data, has true zero. Can't be neg, manipulated by math ops.
 - Ex: weight, speed, age

- Attribute/variable → characteristic / feature of data obj. Vary from one obj. to another.
- Prop. of attributes →
 - Distinctness : =, ≠
 - Order : >, <
 - Addition : +, -
 - Multiplication : *, /
 - Nominal → distinctness
 - Ordinal → distinctness & order
 - Interval → distinctness, order & addition
 - Ratio → all 4.

- Obj, sample → measure / survey members of sample w/o affecting them.
- Controlled study → One grp → some treatment, Other grp + no treatment.
- Web scraping → also in business
 - Use of program/algorithm to extract & process large amt. of data from the web
 - python libraries → BeautifulSoup, Requests, Scrapy
- Web Scraper → program that goes to webpage, downloads contents, extract data out of contents & then saves data to file/db.
- Web crawler → downloading/storing contents of large no. of websites by following webpage

→ Web scraping process →

→ Read resp.

→ Parse and extract

→ Transform data

→ Python libraries →

→ Requests / urllib2 → download webpages

→ BeautifulSoup / LXML → Parsing HTML

→ Pandas / Matplotlib / ~~Numpy~~ Numpy / Scipy → handling / transforming data

→ Pandas / CSV → read CSV files

→ Data Quality →

→ Dirty data → incomplete, noisy/error, inconsistent, intentional, user generated

→ Imp. data quality → better decision making, more confidence, ↓ risk, ↑ results

→ Data cleaning → detecting / correcting corrupt inaccurate data from dataset.

→ Detect, resolve, treat

→ Outliers → data distinctly diff. from other ds,

→ Stat → science of data → collecting, classifying, analysing numerical info, & providing conclusions for it.

→ Descriptive stat, → final result table, graph etc.

→ Inferential stat, → using sample to draw conclusions about pop final result prob/struct

estimating parameters

taking sample stat

& solving subjs abt pop param

hypothesis test

Page No.

sample data need to answer research q's

- Handling missing data →
 - Drop obj → only when we sure missing data is not informative.
 - Drop feature → " "
 - Impute missing → replace with avg / median

- Measures of central tendency → weighted mean
- Mean (\bar{x}) = $\frac{\sum x_i}{n}$, avg value
- Median → middle value
- Mode → most common value

→ Trimmed mean → arrange samples in order & trimming equal no. of them from each end.
 Used for sample size n , p% trimmed mean then no, large, if data pts. to be trimmed → $n p/100$.
 extremely skewed distributions

- Mean → adds more information
- adv, →
 - accounts for all available info, can be combined with other means to find overall mean.
 - easy way to represent entire data with single no.
 - each data set → unique mean value
- 缺点, →
 - very sensitive
 - easily affected by outliers
 - only used with interval/ratio data.

→ Median →

→ adv, →

→ not affected by outliers

→ easily explained as middle value.

→ each set → unique median value.

→ disadv, →

→ value perceived as fit in.

→ cannot be utilized for algebraic treatment

→ Empirical formula: Mode = 3Median - 2Mean

→ Mode →

→ adv, →

→ quick & easy to compute

→ unaffected by outliers

→ used on only level of measurement

→ disadv, →

→ terminal statistic

→ irrelevant if no repeating values

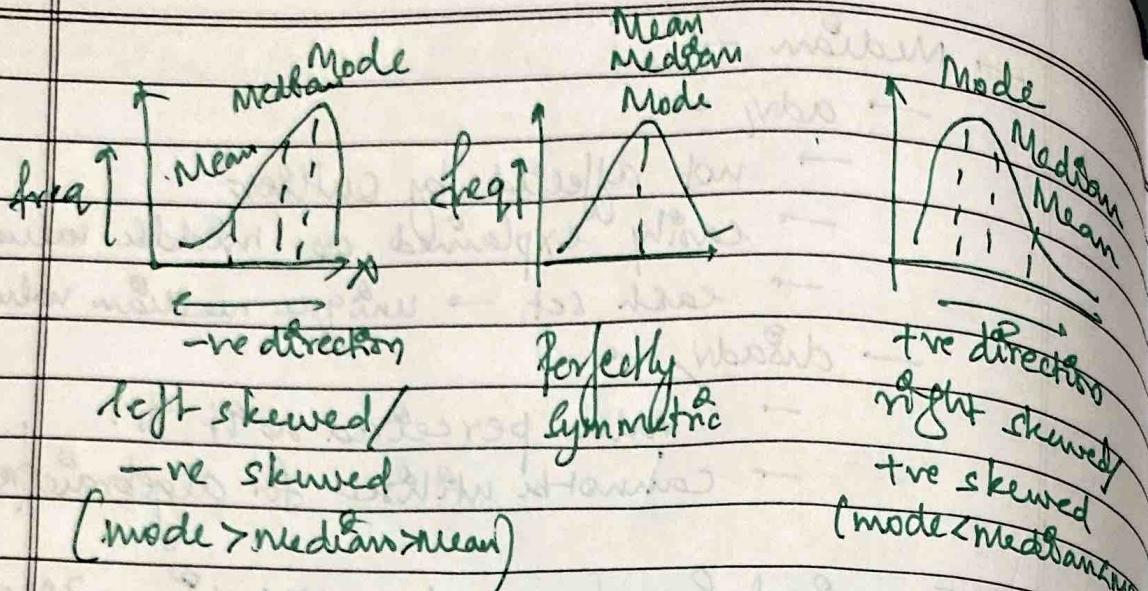
→ many values have same count →

→ Skewness → measure of asymmetry of distribution about its mean.

→ Symmetric distribution → left & right sides equally balanced about mean. Mean, median & mode same.

→ Skewed distribution → "not balanced. Mean & median further toward skew than mode."

↑ distance b/w median & mean → ↑ skewness



Nominal → mode

Ordinal → median

Interval/ratio (not skewed) → mean

" (skewed) → median

Absolute ↑ Relative

→ Measures of dispersion → how homogeneous/heterogeneous data is.

→ Absolute → exp. range in terms of avg. of deviation (the mean/ SD, etc.)

→ Relative → distribution of two/more datasets.

Ex: coeff. of mean, coeff. by SD, etc

→ Range → highest value - lowest value

→ Advantages →

→ easy to calculate

→ simplest of measure of dispersion

→ Ind. of change of origin

→ Disadvantages →

→ affected by fluctuations

→ dep. on change of scale

→ affected by outliers

→ not reliable

→ Percentile → comp. b/w particular value & rest of the values in dataset.

$$\text{Percentile rank } R = \left(\frac{P}{100} \right) * (N+1)$$

→ Quartile → divide list of nos into quarters.

Q_1 → middle no/ b/w smallest no & median
 Q_2 → median
 Q_3 → middle no/ b/w median & largest no

$$Q_1 = 0.25(n+1)$$

$$Q_2 = 0.5(n+1)$$

$$Q_3 = 0.75(n+1)$$

→ Interquartile Range $\rightarrow Q_3 - Q_1$

↓
distance b/w 25th & 75th percentile

→ Variance → measure of spread of recorded values
 ↓
 small → closer to mean
 large → far from mean

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

$$\text{Sample variance } S^2 = \sum (x - \bar{x})^2 / (n-1)$$

$$\text{Population variance } \sigma^2 = \sum (x - \bar{\mu})^2 / N$$

→ Standard Deviation → deviation of elem of dataset from mean value of distribution

→ Histogram →

mode → lower + upper (highest bin height)

$$Q_1, Q_3, \text{median} \rightarrow \text{lower} + \left[\frac{(pos - CF)x}{\text{freq}} \right] \text{class width}$$

$$Q_1 \rightarrow pos = 0.25(n+1), \text{median} \rightarrow pos = 0.5(n+1), Q_3 \rightarrow pos = 0.75(n+1)$$

$$\text{mean} \rightarrow \frac{\sum f(\text{lower} + \frac{\text{upper}}{2})}{N}$$

→ Data visualisation → translating data to visual context.

• Histograms → used for continuous data. area prop, freq., width = class interval

$\sum x_i$
learn histogram construction

$$\text{Bin size} = \frac{\text{IQR}(x)}{3\sqrt{n}}$$

→ Scatter plot → shows whether two vary correlated explanatory var. (ind.) → x-axis & response var. (dep.) → y-axis.

More data cluster along one line → more the variables are correlated.

Pearson coefficient $r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{(n\sum x^2 - (\sum x)^2)(n\sum y^2 - (\sum y)^2)}}$

$x \rightarrow$ value of data point $y \rightarrow$ value of data point

$y \rightarrow$ " on x-axis

$n \rightarrow$ no. of data points

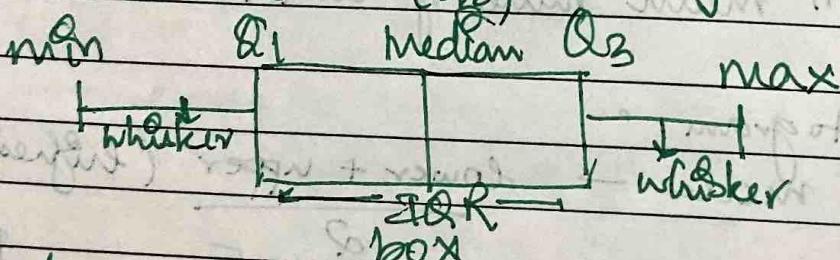
$r > 0 \rightarrow$ +ve correlation

$r < 0 \rightarrow$ -ve r

$r = 0 \rightarrow$ no r

$r = +1/-1 \rightarrow$ perfect correlation

→ Boxplot → way of summarising set of data on interval scale.



Useful when distribution skewed & helps find outliers.

$$\text{min} = Q_1 - 1.5(\text{IQR})$$

$$\text{max} = Q_3 + 1.5(\text{IQR})$$

+ve skew → Q₂ closer to Q₃ than Q₁

-ve skew → Q₂ closer to Q₁ than Q₃

Page No.

→ Bar chart → summarises categorical data.
x-axis → categories, y-axis → counts, perc., etc.

- Heatmap → shows magnitude of phenomenon as color in 2D. Show relation b/w two variables plotted on each of the axes.
Heatmaps used in search engine optimisation

process of improving quality & quantity of web traffic to website/webpage from search engines

- Correlogram → variant of heatmap that replaces variables on the axes with numeric variables

- Dot plot/line plots →
- x-axis → range of values, dots → obs.
 - Dot-height → rep. freq. of obs. values
 - Useful when n is small & many repeated values present.