# Statistics for DSI

## Week 1

Types of data

Descriptive statistics : description & summarisation of data

Inferential stats : draw conclusions from data
- we possibility of chance (probability)

Data — ① Categorical

② Numerical
↳ Discrete
↳ Continuous

Data collection requires one of the following scale of measurement : nominal, ordinal, interval, or ratio.

① 1. Nominal : labels, names

eg : Name, Board, Gender, Blood group etc.
- sometimes numerically coded
- no ordering

② ② Ordinal scale          | Dataset ⟨ Cat ⟨ Nominal
                           |              ⟨ Ordinal
                           |         Num ⟨ Interval
                           |             ⟨ Ratio

⟹ properties of nominal data
- order or rank
eg : customer ratings

③ Interval scale

④ Ratio scale
- All prop of Interval data
- Ratio of 2 vals are meaningful

- all properties of ordinal data          eg : height, weight
- Interval b/n values (fixed units)        age, marks etc
- Always numeric (no absolute zero)        numerical vals
- eg : temperature                         added, subtracted
                                           multiplied divided

# Week 2

Describing categorical data — frequency distributions & relative (count)
frequency (ratio)

## Graphical display

chents of categorical data
- bar chart, Pie chart — relative frequency

### Pie chart
- proportions of category

### Bar chart

Pareto charts : categories in bar chart are
sorted by frequency good for ordinal data variables

---

**Question 2** 5 subjects, total 500 Marks

| | | | |
|---|---|---|---|
| Physics | 35 % | ₹ 175 ✓ | $\frac{125 + 50 + 90}{500}$ |
| chemistry | 25 % | 125 | |
| Biology ✓ | 10% | 50 ✓ | = 63% |
| Maths | 18% | 90 ✓ | |
| Hindi | 12% | 60 ✗ | |

Total Players 200

| Academy | No of Players |
|---------|---------------|
| A | a |
| B | b |
| C | 50 |
| D | d |
| E | 75 |

$a + b + d + 125 = 200$

$a + b + d = 75$

② ④ $\dfrac{a + d + d}{100} = \dfrac{75}{200}$

mode:     bi modal  multi modal

median:     not available unless the data can be put into
            order

total 50 students          $31.5 \times 50$

Grades     A   25    20

           B   32.5  26                  $\dfrac{15}{20}$      $\dfrac{22.5}{25.0}$

           C   22.5 %  %

           D   20 %  %

14   15   20   23   40

_____

Week 3    I measures of Central tendency

Sample mean  $\bar{x} = \dfrac{x_1 + x_2 + \cdots + x_n}{n}$          $n$ — Sample size

                                                                      $N$ — Population size

Population mean  $M = \dfrac{x_1 + x_2 + \cdots + x_N}{N}$

Quiz

| nos | freq |
|-----|------|
| 2   | $x+6$ |
| 6   | $x+2$ |
| 11  | $x-3$ |
| 14  | $x$ |

(1)

$mean = 5.63$

$$\frac{2(x+6) + 6(x+2) + 11(x-3) + 14x}{x+6+x+2+x-3+x} = 5.63$$

$$\frac{2x+12 + 6x+12 + 11x-33 + 14x}{} = 5.63$$

$4x+5$

$$\frac{33x-9}{4x+5} = 5.63$$

$33x - 9 = 5.63(4x+5)$

$= 22.52x + 28.15$

$33x - 22.52x = 28.15 + 9$

$10.48x = 37.15$   $x = 3.54$

$= 4$

Ans 4

① Mean adding constant to each value

new mean $= mean + c$

② Mean multiplied by constant

new mean $= mean * c$

\* medium & mode

① medium

- data Inc order

Sample mean is sensitive to outliers
whereas sample median is not

\* adding constant

New medium = old medium + c

\* Multiplying a constant

new medium = old medium \* c

① If no of obs odd

median $\frac{n+1}{2}$ elem

② If even

mean mean of $\frac{n}{2}$ & $\frac{n}{2} + 1$ obs.

② Mode — most freq

- If no val occurs more than once, no mode

① adding const

new mode = old mode + c

② multiply

new mode = old mode \* c

① measure of dispersion

Range, Variance, std dev.

Range = Max - Min

Range is sensitive to outliers

## Variance

Population variance $\sigma^2 = \dfrac{(x_1 - M)^2 + (x_2 - M)^2 \cdots + (x_N - M)^2}{N}$

Sample variance $s^2 = \dfrac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 \cdots + (x_n - \bar{x})^2}{n-1}$

* ① adding a constant

  new variance = old variance

② Multiplying a constant

  new variance = $c^2$ * old variance

## Standard deviation

$s = \sqrt{\dfrac{(x - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n-1}}$    pop std $= \sqrt{\dfrac{(x - \mu)^2 + \cdots x_{n-2}}{N}}$

sample std $= \sqrt{\text{sample variance}}$

① adding constant : new std = old std
② multiply constant : new std = $c$ * old std

# Percentile, Quantiles & Interquartile range

## Percentiles

1. Arrange the data in the order of position

2. If $np$ is not an integer   smallest int $> np$ position of $100P$.

   Avg values in positions
3. If $np$ is an int. $\frac{V(np) + V(np+1)}{2}$ in the $100P$

④. outliers $< (Q_1 - 1.5 \, IQR)$ and $> (Q_3 + 1.5 \, IQR)$

## Quartiles

mean $= 16$    $\frac{Sum}{10} = 16$

Sum $= 160$

Sum $+ 6 = 166$    $\frac{166}{10} =$

Sample std $= 10$    mean $= 16$

variance $= 100$    $= \frac{(x_1 - \bar{x})^2 + \cdots}{9}$

$Q_1 = \bar{x} | 2 | \cdots = 900$

$= 900 - (10 - 10)^2$

$+ (6 - 10)^2 = \frac{900}{900 + 16} = \frac{916}{9}$

(9) 97, 69, 62, 71, 47    mean @ 18.2

$(2.8)^2 + (91)^2 + (-68)^2 + (75)^2 + (7171)^3$

⑩

⑪  889.48 + 17.64 + 38.74 + 7 57 + 7791 ···779.79

4

(k) 49  38  41  41  96  94  101

38.  40  41  41  96  99  101

n = 7          10th - .1 × 7 = 07 = 1   - 58

40              .5 × 7 = 7.55 87 = 9    909

(l)  IQR = Q₃ - Q₁      75    25

26  73  34  25  77  106  92

i;3   25   26   34   77   92   106

.25 × 7 = 1.75 = 2     Q1 = 25

Q₃ = 75 × 7  5.75 = 6   Q₃  92

outlier  1.5 IQr above orbeli  Q₁ & Q₃
          Q₁ - 1.5 IQr     Q₃ + 1.5 IQR

$Q_1 - 1.5 \text{ IQR}$

$= 25 - 1.5 \times 67$

$Q_3 + 1.5 \text{ IQR} = 92 + 1.5 \cdot 67$

④ 44, 52, 53, 56, 68, 7, 80, 81, 83, 89, 89

90

4

7. 44. 52 53 56 68  80 81 83 89 89

90 $= \dfrac{7 + 12 \cdot 30}{}$

$12 \times \cdot 25 = \dfrac{3}{2}$  4  $= 52 \cdot 5$

$12 \times \cdot 75 = 9$  $= 86$

---

Covariance - <u>Measure of association</u>

| Dev of x | Dev of y | product |
| --- | --- | --- |
| $(x - \bar{x})$ | $(y - \bar{y})$ | $(x - \bar{x})(y - \bar{y})$ |

Population covariance

$$\text{cov}(x, y) = \dfrac{\sum_{i=1}^{N}(x_i - \bar{x})(y_i - \bar{y})}{N}$$

Sample covariance

$$\text{cov}(x, y) = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

## Correlation

Measure of linear correlation b/w 2 numerical
variables

$$r = \boxed{\frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}} = \frac{Cov(x,y)}{\sigma_x \cdot \sigma_y}$$

$$\frac{Cov(x,y)}{S.std_x \times S.std_y}$$

Price = 7.77 × Size + 130

$R^2 = 0.022$     $r = \sqrt{0.022} \cdot 0.149$     $R^2 = 1 - \frac{RSS}{TSS}$

↓ goodness of the fit

$y = mx + c$     $m = slope$
          $c = Interrupt$

$0 \le R^2 \le 1$

correlation ↗ $x$

correlation r

| x | y | $(x - \bar{x})$ | $(y - \bar{y})$ | $(x - \bar{x})^2$ | $(y - \bar{y})^2$ | $(x-\bar{x})(y-\bar{y})$ |
|---|---|---|---|---|---|---|
| 355 | 300 | -106.9 | -74 | 11427.61 | 5476 | 7910.6 |
| 487 | 340 | 25.1 | -34 | 630.01 | 1156 | -853.6 |
| 526 | 400 | 64.1 | 26 | 4108.81 | 676 | 1666.6 |
| 590 | 450 | 128.1 | 76 | 16409.61 | 5776 | 9735.6 |
| 428 | 300 | -33.9 | -74 | 1149.21 | 5476 | 2508.6 |
| 398 | 325 | -63.9 | -49 | 4083.21 | 2401 | 3131.1 |
| 555 | 450 | 93.1 | 76 | 8649 | 5776 | 7075.6 |
| 320 | 400 | -141.9 | 26 | 20135.61 | 676 | -3689.4 |
| 450 | 375 | -11.9 | 1 | 141.61 | 1 | -11.9 |
| 510 | 400 | 48.1 | 26 | 9623.61 | 676 | 1250.6 |

mean = 461.9   374

76358.29   28090

$$\frac{28723.8}{\sqrt{76358.29} \times \sqrt{28090}}$$

$$= \frac{28723.8}{276.33 \times 167.60}$$

$$= \frac{28723.8}{46312.908} = 0.6202 1$$

Revision

Week 2

| | | five | 150 | 150 |
|---|---|---|---|---|
| W1 | 1+1 | 2 | r 175 | real |
| IND) 1+1 | | 2 | 1 125 | 125 |
| AM1 1+1+1+1+1 | | 5 | 150 | 100 |
| 7AK 1 | | 1 | 300 | |
| SL 1 | | 1 | 1000 | |
| CNG 1 | | 1 | | |

$\underline{\phantom{12}}$ 12

(2) total no of tally

$150 + x + y + 250 + 300$

$= 700 + x + y$

$\dfrac{250}{700 + x + y} = 0.25$    $\dfrac{300}{700 + x + y} = 0.3$

$250 = 0.25(700 + x + y)$    $300 = 0.3(700 + x + y)$

$1000 = 700 + x + y$    $x + (x - 300) = 300$

$x + y = 300$    $2x - 300 = 300$

$y = (x = 380)$    $y = 300 - 175$

$\dfrac{x}{700 + 300} = 0.175$    $= 125$

$x = 0.175(1000) = 175$