

Descriptive Inferential Statistics

DA
① - read

② - analysis

Statistical methods

ML
① - read data

② - clean data

③ - preprocess data

④ - model building

70% → DA
30% → ML
get trigger

Hand-drawn diagram illustrating the data science process:

- Data (past)** (circled) leads to **Analysis** (boxed).
- Analysis** leads to **Data (future)** (boxed).
- Data (future)** leads to **pred** (written).
- pred** leads back to **Data (past)**.
- A **m** (model, circled) is shown between **Analysis** and **pred**.

↓
- Analysis
- patterns

↓ prediction (ML) | Hypothesis
↓
Intervention

- Measures of Central Tendency
- Measures of Dispersion / Deviation
- Measures of Shape
- measures of relation

$$\frac{10 + 11 + 12 + 13}{4} = \frac{46}{4} = 11.5$$

① mean :- Sum of all obsⁿ / no of the obsⁿ (\bar{x}) / (μ) $\bar{x} = \frac{\sum x}{n}$

② median :- middle value / which is divided into two equal parts of data

③ mode

median

$$\text{mean} = \frac{7+15+13+6+4+2+10}{7} = 57.14 = \underline{8.14}$$

median =

median:- $2, 4, 6, \textcircled{7}, 10, 13, 15 \Rightarrow m = 7$

$$\begin{array}{c|c} 50.1 & 50.1 \\ \hline n \end{array}$$

$$n = 8$$

2, 4, 6, 7, 10, 12, 15, 100

$$\text{mean} = 15718 = 19.6$$

$$\text{median} = \frac{n^{\text{th}} + (\frac{n}{2} + 1)^{\text{th}}}{2}$$

$\text{median} = 8.5$
 $\text{median} = 8$
 $9.6 \geq 8.5$

mean = 19.6

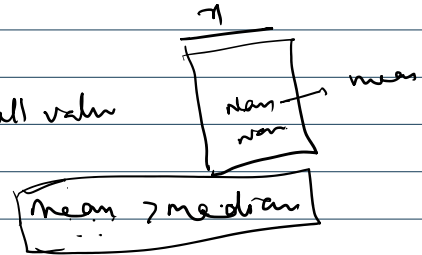
mean = 8.14

mean = 8.5

$$\frac{19.6 + 8.5}{2} = \frac{1 + 5}{2} = \frac{7 + 10}{2} = \frac{17}{2} = 8.5$$

Note - mean is effected by the out lias but not median

- CT - ① Imputation of the data null value
② Identifying the distn



mode - High concentration / highly repeated value

2, 4, 5, 1, 2, 4, 1, 2, 2 → mode: 2 - 4 times
1 - 2 times

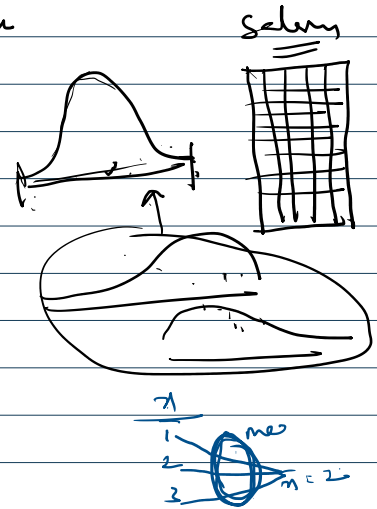
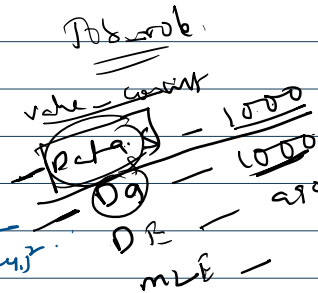
Numerical

- ① Imputation of null value where mean / median

Categorical

- ① Imputation of null value we use mode

②



Measures of Dispersion

H - L

- ① Range: - Max - min

- ② SD - Standard Deviation: Distance the samples

- ③ variance

SD - Distance the original point / observation to mean

$$SD = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} = \sigma \text{ (sigma) (Root mean squared error)}$$

variance - $\sigma^2 = \frac{\sum (x - \bar{x})^2}{n}$ (mean squared error)

SD is also identifies of the distn

Diagram showing a box plot with 'sum' and '10 mins'.

x	(x - \bar{x})	(x - \bar{x}) ²
2	2 - 3.6 = -1.6	2.56
4	4 - 3.6 = 0.4	0.16
6	6 - 3.6 = 2.4	5.76

$\sqrt{\frac{8.48}{5}} = 1.3$

Diagram showing a box plot with '2H', '4H', '6H', '4H' and 'mean'.

$$\frac{2 \times 4 + 6 \times 4 + 4 \times 2}{5} = 16.6 = 1.6$$

SD = 1.4

0.1

0.2

4	$4 - 3.6 = 0.4$	0.16	$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$	6+	$\left. \begin{array}{l} 6+ \\ 4+ \\ 2+ \end{array} \right\}$	0.1
6	$6 - 3.6 = 2.4$	4.86		4+		0.2
4	$4 - 3.6 = 0.4$	0.16		2+		
2	$2 - 3.6 = -1.6$	2.56				
		<u>10.30</u>				

$$\sigma = \sqrt{\frac{10.3}{5}} = \sqrt{2.06}$$

$$\frac{16}{5} = 3.2$$

$$SD = 1.4 \quad \xrightarrow{2 \text{ mean} = 3.6} \quad \boxed{2.2}$$

x	$(x - \bar{x})$	$(x - \bar{x})^2$	$(x - \bar{x})^3$
2	$2 - 3.2 = -1.2$	1.44	
4	$4 - 3.2 = 0.8$	0.64	
4	$4 - 3.2 = 0.8$	0.64	
4	$4 - 3.2 = 0.8$	0.64	
2	$2 - 3.2 = -1.2$	1.44	
		<u>4.80</u>	

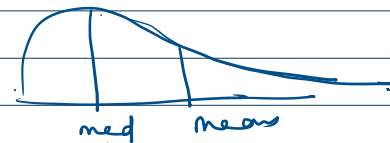
$$\sigma = \sqrt{\frac{4.80}{5}} = \sqrt{0.96}$$

$$SD = 0.97 \quad \uparrow \quad 3.2$$

`diamond_df['price'].describe()`

```
count    53940.000000
mean     3932.799722
std      3989.439738
min       326.000000
25%      950.000000
50%     2401.000000
75%     5324.250000
max     18823.000000
Name: price, dtype: float64
```

mean > med positive (skewed)



SD > mean

