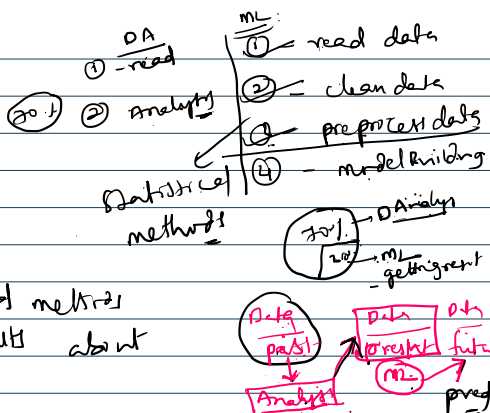


Statistical methods
Descriptive Inferential Statistics

Stats - Representation of your data
in terms of graph or tables
and applying some statistical methods
to get interpretations / results about
your data



Descriptive Stat - As Analysis of past data and
get the pattern for past data which are useful
for present / future data.
prediction (ML) | hypothesis
↓
inferential

Descriptive :-

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

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

Central Tendency:

- ① mean :- Sum of all observ / no. of the observ $(\bar{x}) / (N)$ $\bar{x} = \frac{\sum x}{n}$
- ② median :- middle value / which is divided into two equal parts of data
- ③ mode

① 7, 15, 13, 6, 4, 2, 10, 100

$$\text{mean} = \frac{7+15+13+6+4+2+10+100}{8} = 19.6$$

$$\text{median} = \frac{7+10}{2} = 8.5$$

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

② 7, 15, 13, 6, 4, 2, 10, 100

$$\text{mean} = 19.6$$

mean = 19.6
mean = 8.5

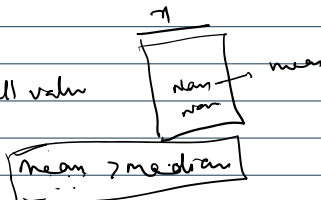
median = 7.5
median = 7

$$19.6 > 7.5$$

$$= \frac{7+10}{2} = 8.5$$

Note - mean is affected by the outliers
but not median

- CI - ① Imputation of the data null value
② Identifying the data



mode - High concentration / highly repeated value

2, 4, 5, 1, 2, 4, 1, 2, 2 \rightarrow mode: 2 - 4 times

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

Numerical

① Implication of null value never mean / median

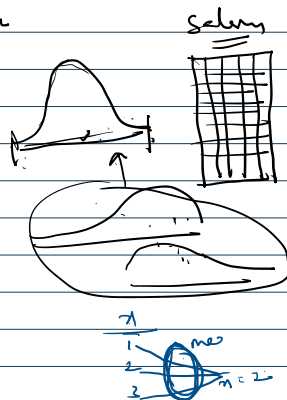
Categorical

① Implication of null value we use mode

②



Distance
value = constant
1000
1000
DE = 290
mbe -



Measures of Dispersion

H - L

- ① Range: Max - min
- ② SD - Standard Deviation: Distance b/w samples
- ③ Variance

SD: Distance b/w original point / observation to mean

$$SD = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} = 0 \text{ (same)} \text{ (Root mean squared error)}$$

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

*** SD is also identifier of the distn



2ft 2ft 2ft

$$\frac{2+4+6+4+2}{5} = 16/5 = 3.2 \text{ mean}$$

n	(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
4	4 - 3.6 = 0.4	0.16
2	2 - 3.6 = -1.6	2.56
		10.20

$$\sqrt{\frac{10.2}{5}} = \sqrt{2.04}$$

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$\sigma = \sqrt{\frac{10.2}{5}} = \sqrt{2.04}$$

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

SD = 1.4

$$\sigma = 0.4$$

n	(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	
		4.80	

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

SD = 0.96

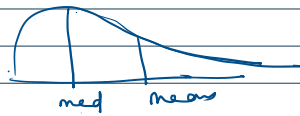
$$\sigma = 0.96$$

```

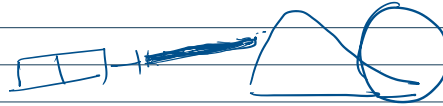
diamond_df['price'].describe()
count    53940.000000
mean     3932.799722
std      3989.439738
min      326.000000
25%     950.000000
50%     2401.000000
75%     5324.250000
max     18823.000000
Name: price, dtype: float64

```

mean > med positive skewness

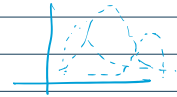


SD > mean



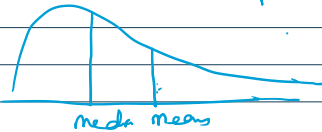
Measures of shape:

Skewness → scatterness



① mean > median → positive skewness

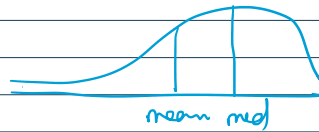
$Sk > 0$



$Sk > 0.5$

② mean < median → negative skewness

$Sk < 0$

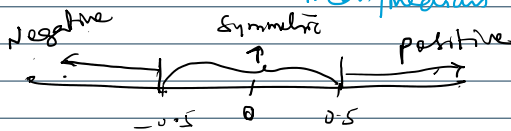


$Sk < -0.5$

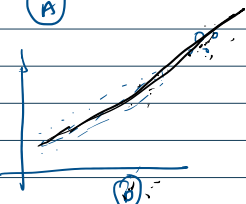
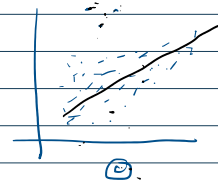
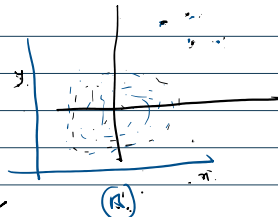
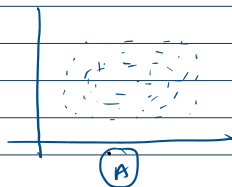
③ mean = median → symmetric skewness / Normal curve

$Sk = 0$

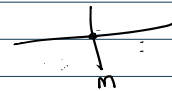
$-0.5 \leq Sk \leq 0.5$



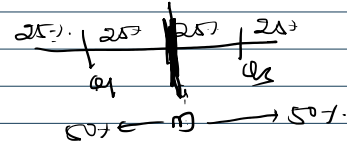
Outliers:-



Correlation:



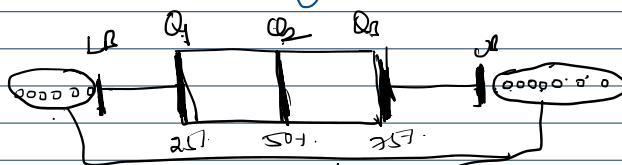
Quantiles:



Q1 - Lower Quantile - 25%

Q3 - Upper Quantile - 75%

Q2 - Median - 50%



Outliers

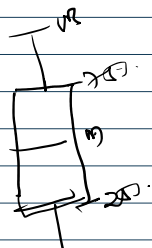
ign.

Inter-Quartile Range:

$UR = Q_3 - (1.5 \times IQR)$

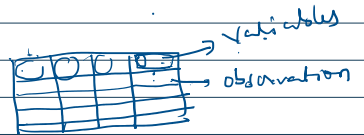
$LR = Q_1 - (1.5 \times IQR)$

... $Q_0 = Q_1$



Variables

$UR = 42.11$
 $LR = a_1 - (1.5 \times 5.47)$
 $WR = a_2 - a_1$



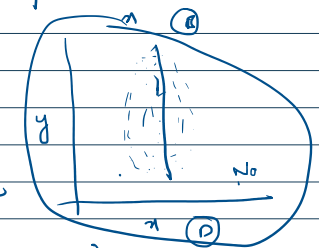
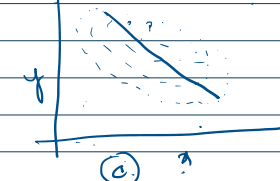
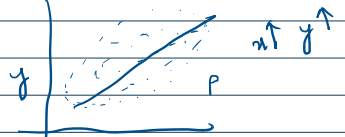
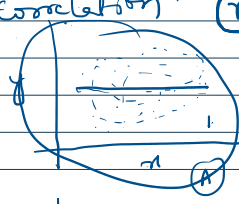
Measures of relation:

relation b/w two variables/columns is called "correlation"

① positive correlation

② negative correlation

change in one variable leads to change in other variable is called having relation



① If we have opposite direction in the change of variable is called negative

$x \uparrow y \downarrow$ | $x \downarrow y \uparrow$

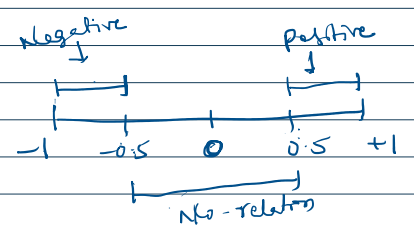
$r < 0$

② If we have same direction in the change of variable is called positive correlation

$x \uparrow y \uparrow$ | $x \downarrow y \downarrow$

$r > 0$

$r \approx 0$ No correlation



① Because if the age decreases my salary is also decreasing increasing in salary is also increasing

$Corr = 0.7$
 70%

$age \downarrow$ salary \downarrow
 $age \uparrow$ salary \uparrow

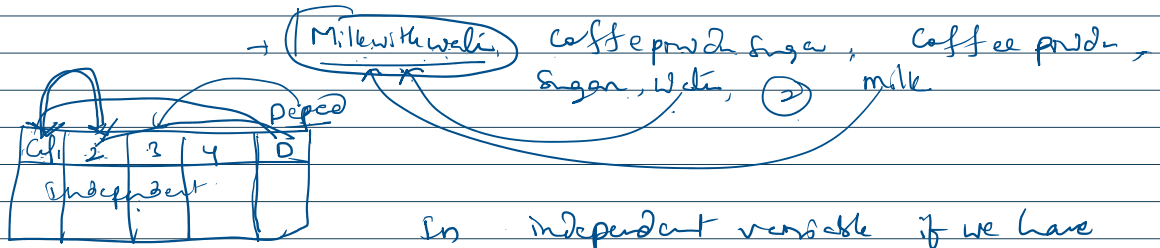
positive $0.5 < Corr < 1$

-0.6
 60%

$age \uparrow$ salary \downarrow
 $age \downarrow$ salary \uparrow

negative $-0.5 < Corr < -1$

Coffee \rightarrow (Milk), (Water), Coffee powder, Sugar ①



In independent variable if we have relation b/w the variable is called "multi-collinearity"

note - def remaining cell are independent

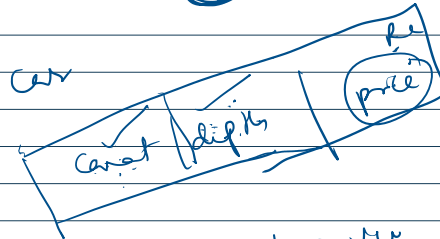
	1	2	3	4	5	6
1	1	0.028	0.18	0.92	0.98	0.95
2	0.028	1	-0.3	-0.011	-0.025	-0.029
3	0.18	-0.3	1	0.13	0.2	0.18
4	0.92	-0.011	0.13	1	0.88	0.87
5	0.98	-0.025	0.2	0.88	1	0.86
6	0.95	-0.029	0.18	0.87	0.86	1

To identify the multicollinearity we have $\sqrt{R^2}$ Variance Inflation factor $VIF < 10$

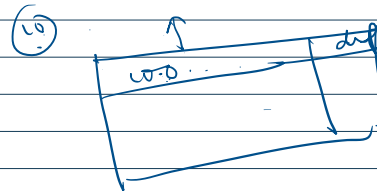
Heatmap visualization showing the correlation coefficients between variables in the 'diamonds' dataset. The color scale ranges from -0.2 (dark purple) to 0.6 (yellow). The 'price' variable is circled in blue, and the 'depth' variable is highlighted with a red vertical line.

	carat	depth	table	price	x	y	z
carat	1.00	0.18	0.13	0.02	0.18	0.15	
depth	0.18	1.00	0.13	0.01	0.88	0.87	0.86
table	0.13	0.13	1.00	0.88	0.97	0.97	
price	0.02	0.01	0.88	1.00	0.97	0.95	
x	0.18	0.88	0.97	0.97	1.00	0.95	
y	0.15	0.87	0.97	0.95	0.95	1.00	
z		0.86	0.97	0.95	0.95	0.95	1.00

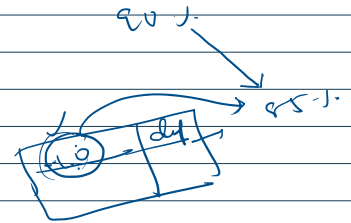
cost
 4
 dep. th.
 price → dep



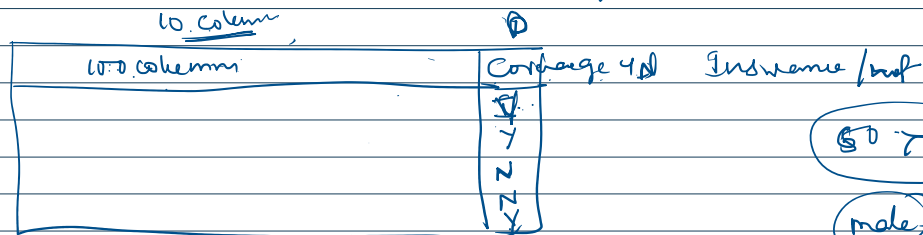
गु.त.



100 mm



{ patient demography details — 10 — ID, MEN, age, gender, employment stat
 patient cancer details — paper, diagnosis, stage, dx, language, surgery
 patient ~~disease~~ details — area, rx, md, visit, F/M
 90 columns



50-2 m/s

male → Sun

- ① age ✓
② gender ✓
③ employment ✓
④ population ✓
⑤ are ✓
⑥ cities ✓

- correlated
- multi

swrmi