

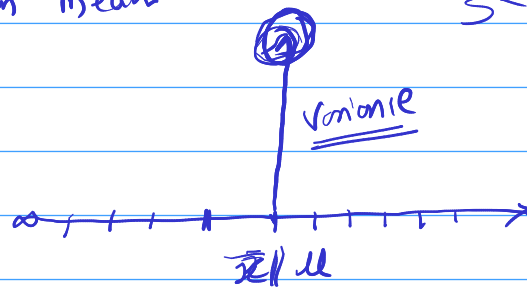
## Standard Deviation

$\bar{x}$  - Sample mean  
&  $\mu$  → Population mean

$$\text{Variance } \sigma^2 = \frac{\sum (x - \mu)^2}{N}$$

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

Variance -

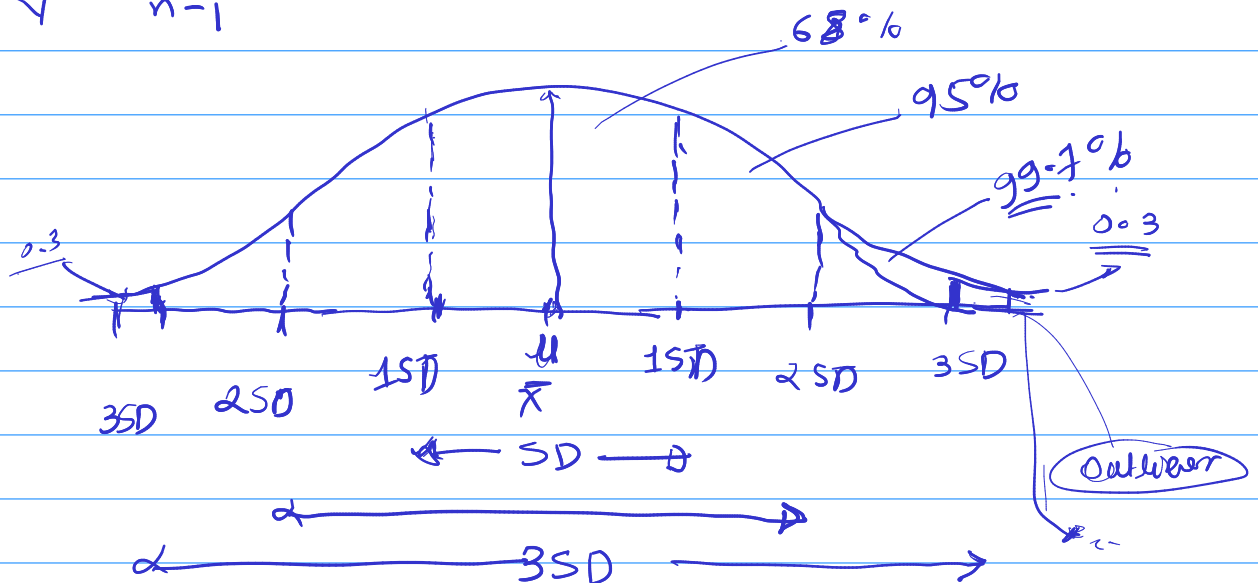


Standard Deviation → is the square root of variance.

It is widely used measure of dispersion that is useful in describing the shape of distribution

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}} = \sigma = \sqrt{s^2}$$

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} \quad s = \sqrt{s^2}$$



Empirical Rule of the data ⇒

### ② Coefficient of Variation

$$\text{Coef of var} = \frac{\text{SD}}{\text{mean}} \times 100$$

① Variance

② SD

③ Coef. of Variation



Ratio of standard deviation to the mean expressed as a percentage

Univariate Analysis about the data

- univariate (single feature)
- Bivariate (2 features)
- Multivariate (multiple feature)

$f_1$   
Categorical

$f_2$   
Numerical

Analyze the (EDA) = Exploratory Data Analysis → data

1) Categorical Data → Frequency Distribution Table

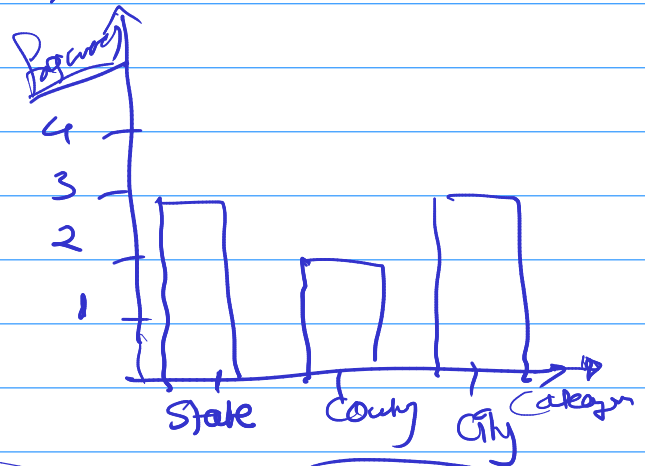
$f_1$

State
City
Country
State
State
City
Country
City

Frequency Distribution

	Frequency
State	3
Country	2
City	3

Bar chart  
Bar graph



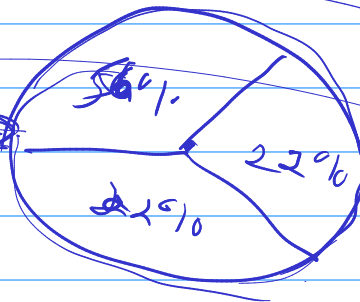
Bar - freq →

For Pie-chart

Calculate Relative freq

freq	Relative freq
3	56.25%
2	22.22%
2	22.22%
	100%

Relative freq



Cumulative Frequency

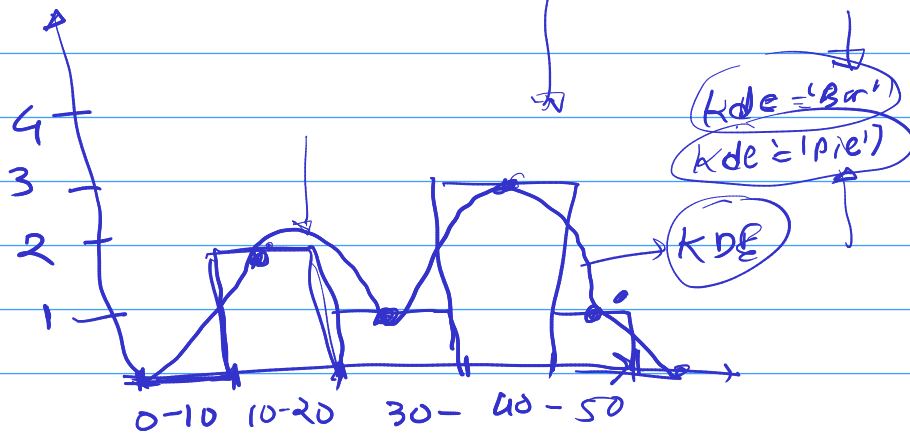
3
5
7

## 2) Numerical - Feature

Age

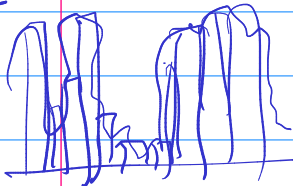
Age	freq
32	0
28	2
35	1
39	3
41	1
12	
15	

## Bar chart - Histogram

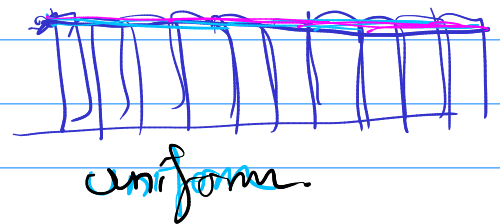
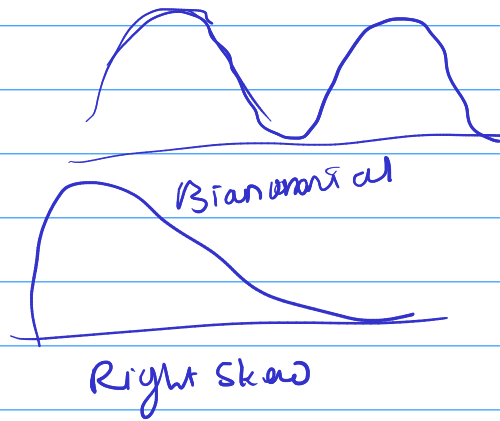


KDE = Kernel Density Function Estimator

Bell curve  
Symmetric curve  
Normal dist  
= Gaussian dist



No-pattern graph



uniform

① → Bar

② - Pie

③ Count plot

④ Distplot → warning

⑤ KDEplot (new)

Boxplot

# Graph for Bivariate Analysis (2 feature) → category & Numerical

Categorical & Categorical

Numerical & Numerical

Numerical & Categorical

## ① Categorical & Categorical

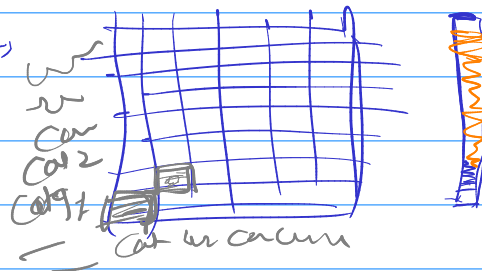
Tip

### ⇒ Contingency Table / Crosstab

<u>Survived</u>		<u>P-class</u>	→	Survive	1	2	3
not survive → 0	1	1		0	42	31	63
Survive → 1	2	2		1	71	118	13
	3	3					

## ② Heatmap

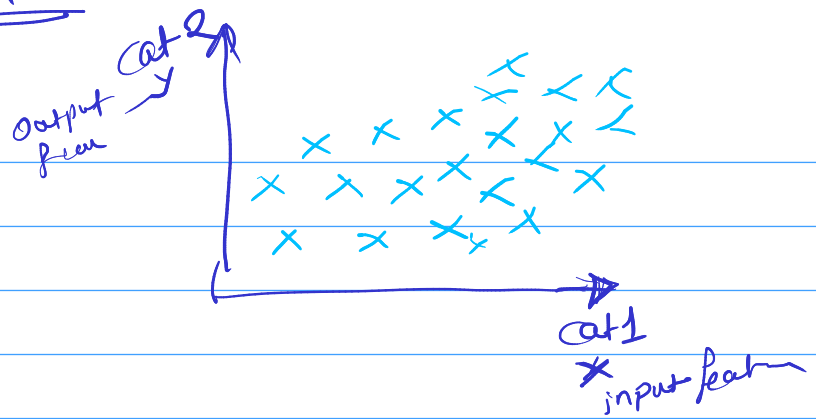
Categorical



## 2) Numerical - Numerical

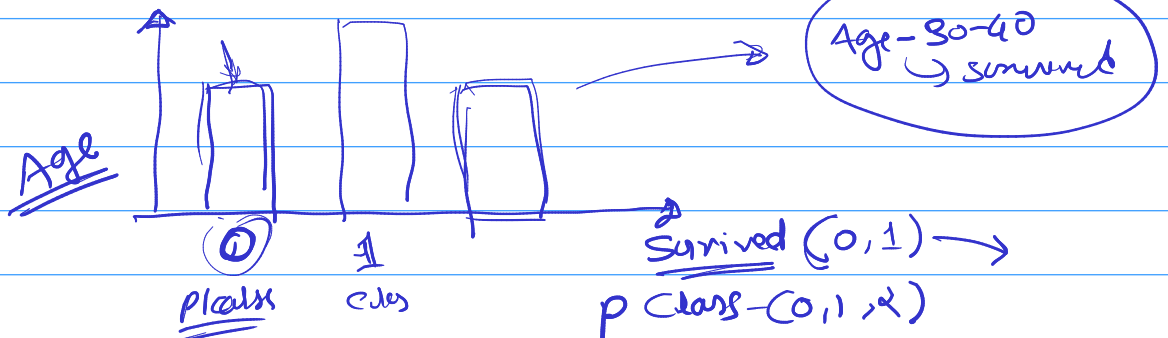
17

Scatter plot

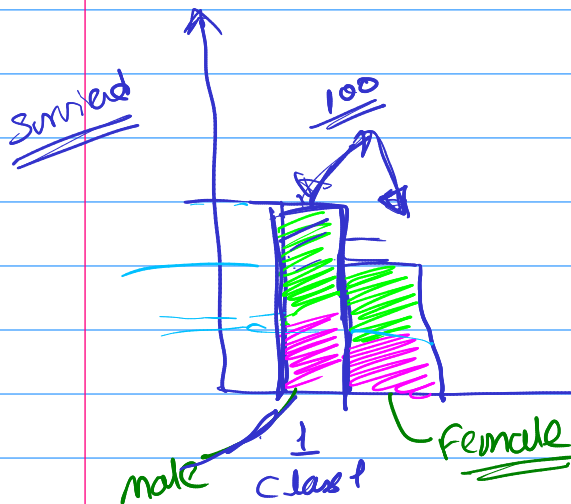


## 3) Categorical - Numerical / Numerical & Categorical.

① Bar



② Hue (Pie - Bar)  
(multivariate Analysis)



pink - Survived  
green - not survived

p class 1 = total (out of 100)

male group is more in class 1

female is less

in male green is not survived

p class pink is survived

in female -

green is not survived

& pink is survived.

multivariate analysis

{ 'sex' } { 'p class' } { 'Survived' }

Male

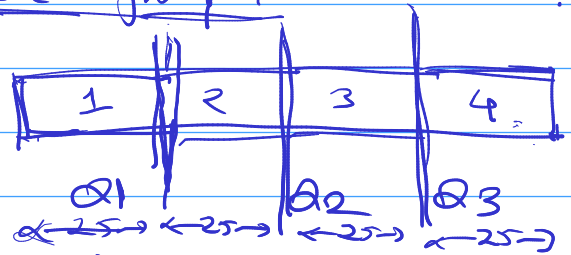
||  
x axis

||  
y axis

# Quantile and percentile

Quantile are Statistical measure use to divide a set of numerical data into equal sized group

Type of Quantile



① Quartile → Divide the data into four equal part  
 $Q_1 = 25^{\text{th}}$  percent     $Q_2 = 50^{\text{th}}$      $Q_3 = 75^{\text{th}}$

② Deciles — we divide data into 10 equal part.

③ Percentiles — Divide the data into 100 equal part

④ Quintiles — Divide into 5 equal part

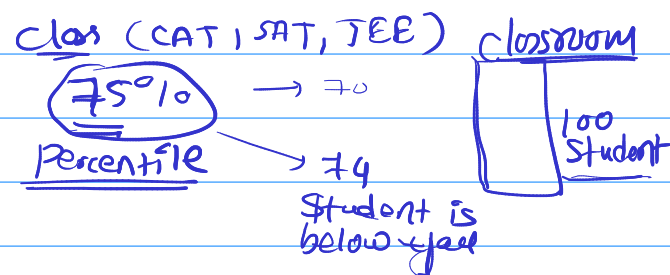
- ① Sort the data —
- ② You get an index location

① Percentile →

$$\text{percentage} = \frac{\text{The no. we got} \times 100}{\text{total number}}$$

$$= \frac{560}{800} \times 100 = 70.75\%$$

A percentile is a statistical measure that represent the percentage of observations in dataset that fall below a particular value.



Formula for

$$\text{Percentile} = \frac{P}{100} (N+1)$$

$P$  → the percentile rank  
 $N$  → total no. of observation.

ex = 78, 82, 84, 88, 91, 93, 94, 96, 98, 99 → sorted

$\begin{matrix} \textcircled{1} & \textcircled{3} & \textcircled{2} & \textcircled{1} \\ \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & 2 & 3 & 4 \end{matrix}$ 
 $\begin{matrix} \textcircled{5} & \textcircled{6} & \textcircled{7} & \textcircled{8} & \textcircled{9} & \textcircled{10} \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 5 & 6 & 7 & 8 & 9 & 10 \end{matrix}$

if you find 75<sup>th</sup> percentile

$$PL = \frac{P}{100}(N+1) = \frac{75}{100}(10+1) = 0.75(11) = 8.25$$

↓  
location

75<sup>th</sup> percentile of data is 97  
= 75<sup>th</sup>

$\frac{8^{th} + 9^{th}}{2} = \text{average} = \frac{96 + 98}{2} = 97$

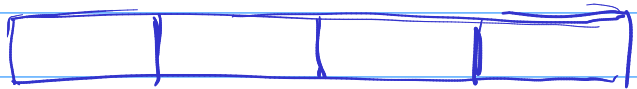
rank of 91 →

$$\text{Percentile rank} = \frac{x + 0.5}{N}$$

$x$  → number of value below the given value  
 $N$  → total no. of dataset

$$\text{rank of 91} = \frac{4 + 0.5}{10} = 0.45 = 45\%$$

5-Number Summary — Aim → to observe the number



① minimum value ( $Q_1 - 1.5(IQR)$ )

② first Quartile ( $Q_1$ ) (25%)

③ median ( $Q_2$ ) (50%)

④ Third Quartile ( $Q_3$ ) (75%)

⑤ maximum value ( $Q_3 + 1.5(IQR)$ )

$$IQR \text{ InterQuartile Range} = (Q_3 - Q_1)$$

6, 213, 241, 260, 281, 290, 314, 321, 350, 1500

$\begin{matrix} \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \end{matrix}$

$$Q_2 = \frac{5.5}{100}(10+1) = 11 = \frac{5.5}{\text{location}} = \frac{5^{th} + 6^{th}}{2} = \frac{281 + 290}{2} = 285.5 \text{ (} Q_2 \text{)}$$

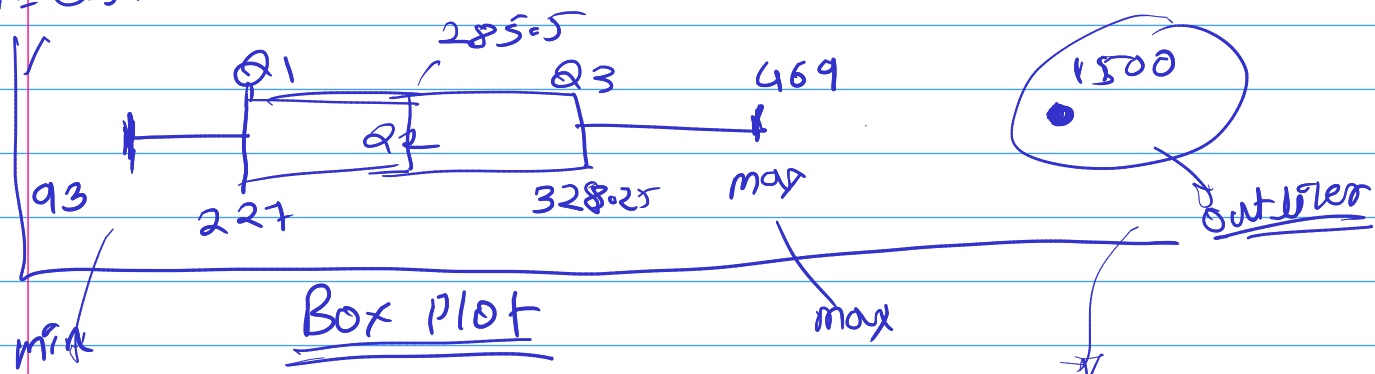
$$Q_1 = \frac{2.5}{100}(11) = 2.75 = \frac{2^{th} + 3^{rd}}{2} = \frac{213 + 241}{2} = 227 \text{ (} Q_1 \text{)}$$

$$Q_3 = 328.25$$

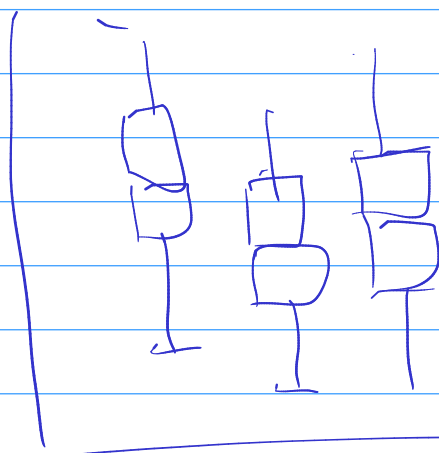
$$IQR = Q_3 - Q_1 = 94$$

$$\min = Q_1 - 1.5 \cdot IQR = 93$$

$$\max = Q_3 + 1.5 \cdot IQR = 469$$



multivariate analysis



use to compare  
2 category also  
using multivariate  
Box Plot analysis.