

Unit-4 - modelling Data Distribution

PERCENTILES

% of data that is ^{at or} below the amount in question.

2 score ÷

how many σ away from μ .

$$2\text{score}(x) = \frac{x - \mu}{\sigma}$$

- if $2\text{-score} > 0$, data point is above avg
- if $2\text{-score} < 0$, data point is below avg
- if 2-score close to 0, data point close to avg.

→ Effects of Linear transformation :

• Adding a constant ($+ \alpha$)

$$\text{mean} = \text{mean} + \alpha$$

$$\text{std-dev} = \text{same}$$

$$\text{median} = \text{median} + \alpha$$

$$\text{IQR} = \text{same}$$

• Multiply a constant ($\times \alpha$)

$$\text{mean} = \text{mean} \times \alpha$$

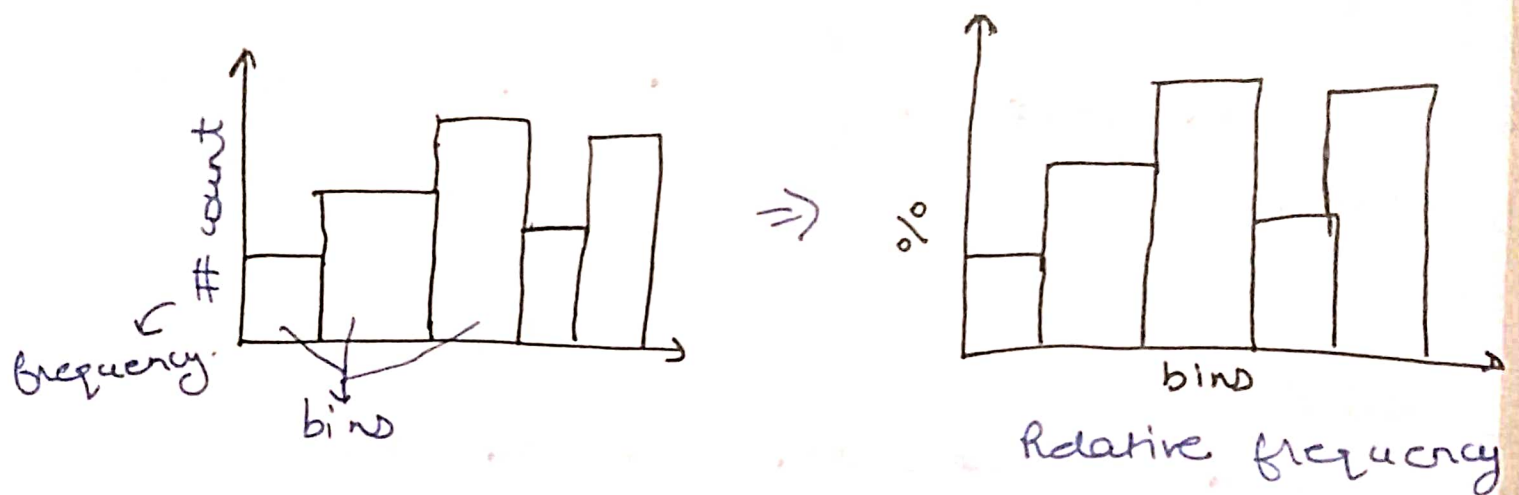
$$\text{std-dev} = \text{std-dev} \times \alpha$$

$$\text{median} = \text{median} \times \alpha$$

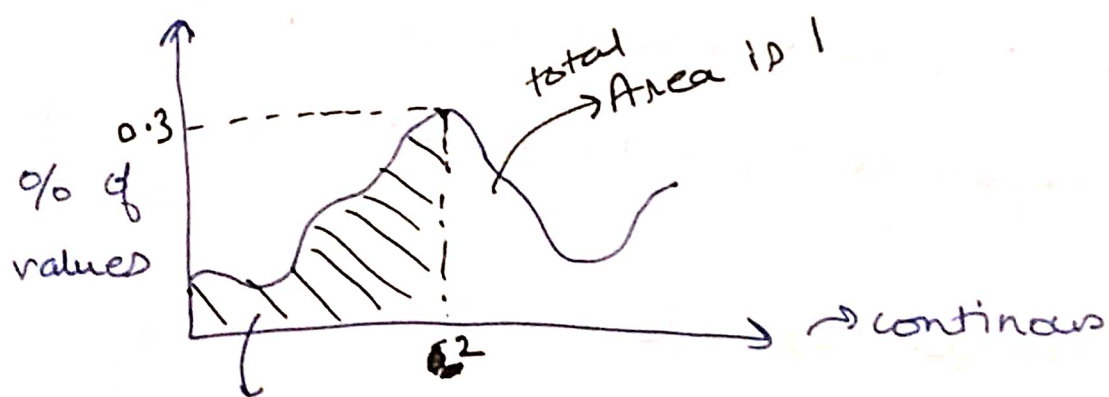
$$\text{IQR} = \text{IQR} \times \alpha$$

⇒ Density curves ⇒

So far we are using histograms to visualize distribution



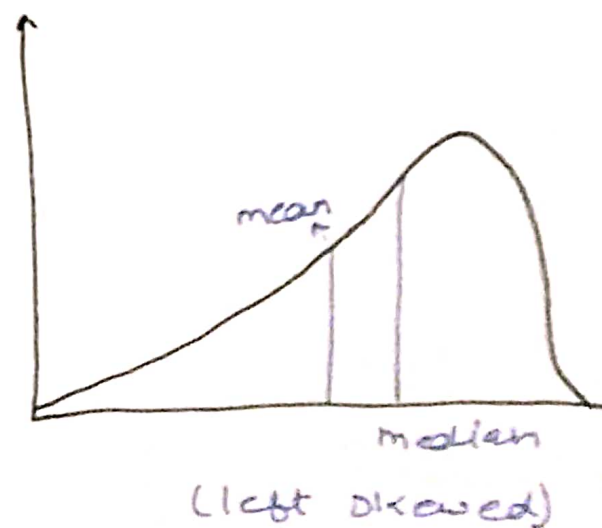
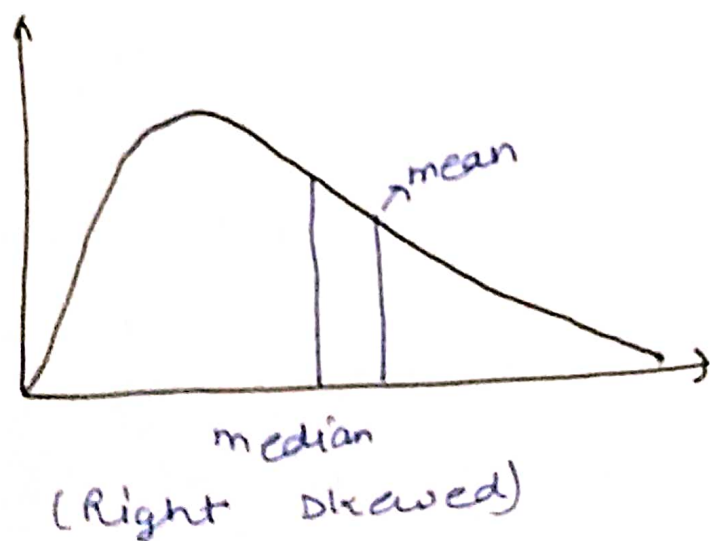
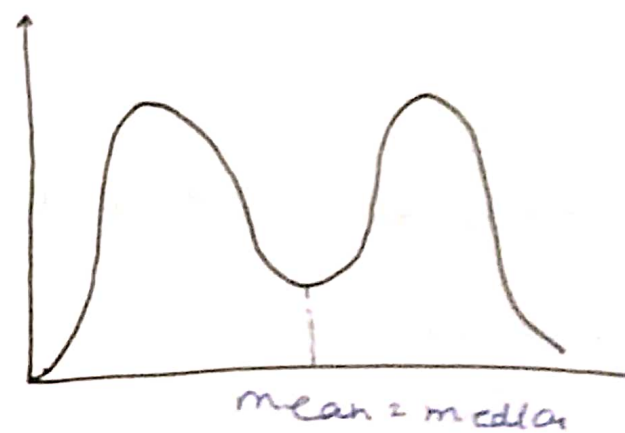
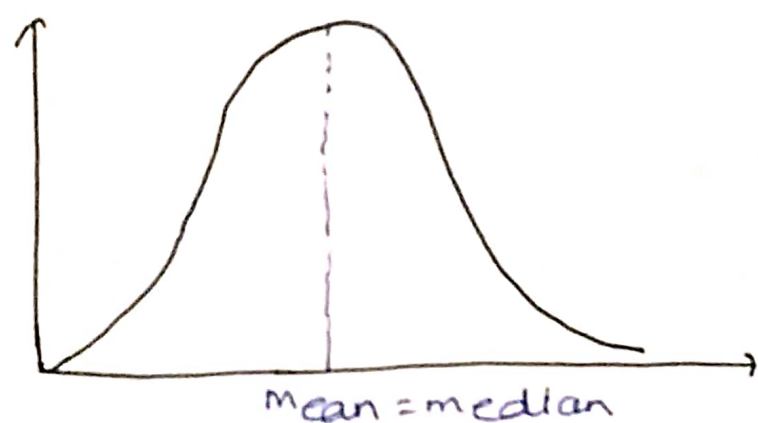
- In density graph we make infinite number of bins with approximately 0 width.



Area = % of values having value < 2.

mean, median, skew from Density curves

- median \rightarrow Point which divides the density curve into two equal areas
- mean \rightarrow take each point, multiply by their frequency and add them up.
 - Balancing point.
- For symmetrical distribution, mean = median.



Normal Distribution and the Empirical Rule

68 - 95 - 99.7 rule

68% within 1 std. dev away from mean



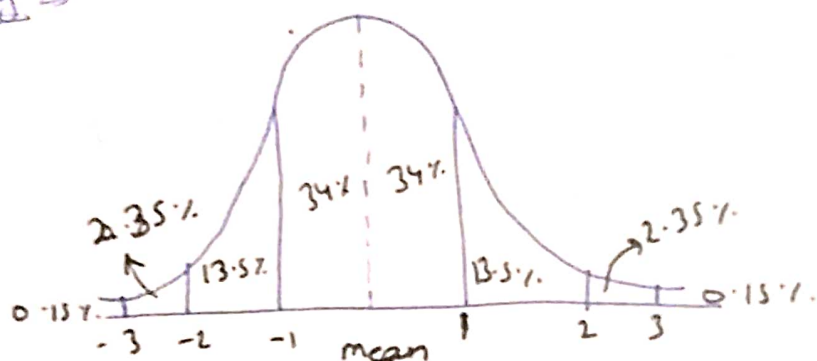
95% within 2 std dev away from mean

99.7% within 3 std dev away from mean

Standard normal distribution \rightarrow

Normal dist. where mean = 0, std dev = 1

\Rightarrow Summary \Rightarrow



\Rightarrow Normal distribution calculation

- mean & σ is given. Find the area of the shaded region



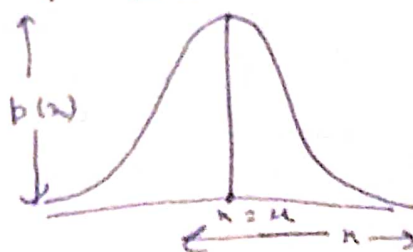
Steps \rightarrow

Step 1 \div Find z score

Step 2 \div look into the table to find corresponding area.

\Rightarrow more on Normal distribution \Rightarrow

$$pdf \Rightarrow p(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$



• Note:-

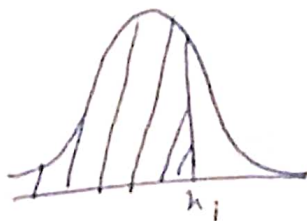
- at a given x , $p(x)$ gives the height of the function, and not the probability.
- For a given number, let's say 5, $p(5) \geq 0$.
(prob. is area under the graph), as in this case the area $= 0$.
- To find probability we have to give some range and integrate.

$$\int_{4.9}^{5.1} p(x) dx$$

\Rightarrow Cdf - cumulative density function.

$$cdf(x) = \int_{-\infty}^x p(x) dx$$

- For a given x_i , $cdf(x)$ gives the probability the $x < x_i$,
i.e. the area under graph \downarrow



To find the prob of x lying in between x_1 & x_2 .
 $p(x_1 < x < x_2)$ it will be,
 $cdf(x_2) - cdf(x_1)$.

