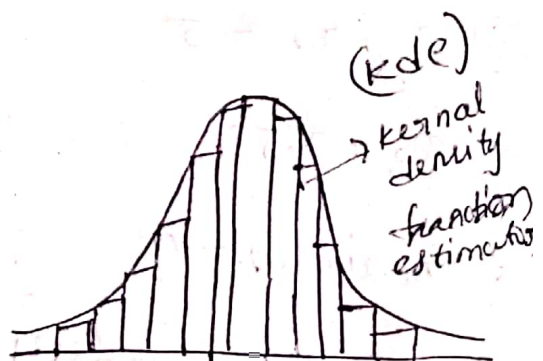
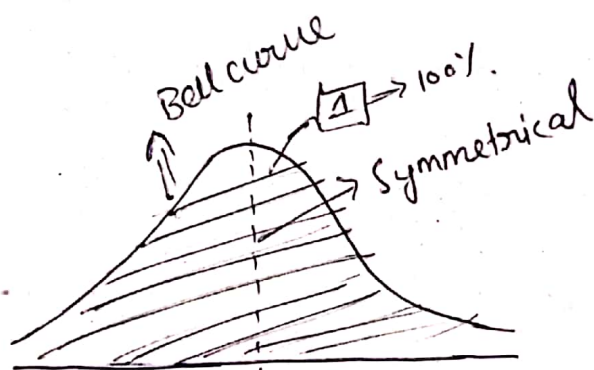


## Statistics part-3

### Agenda!

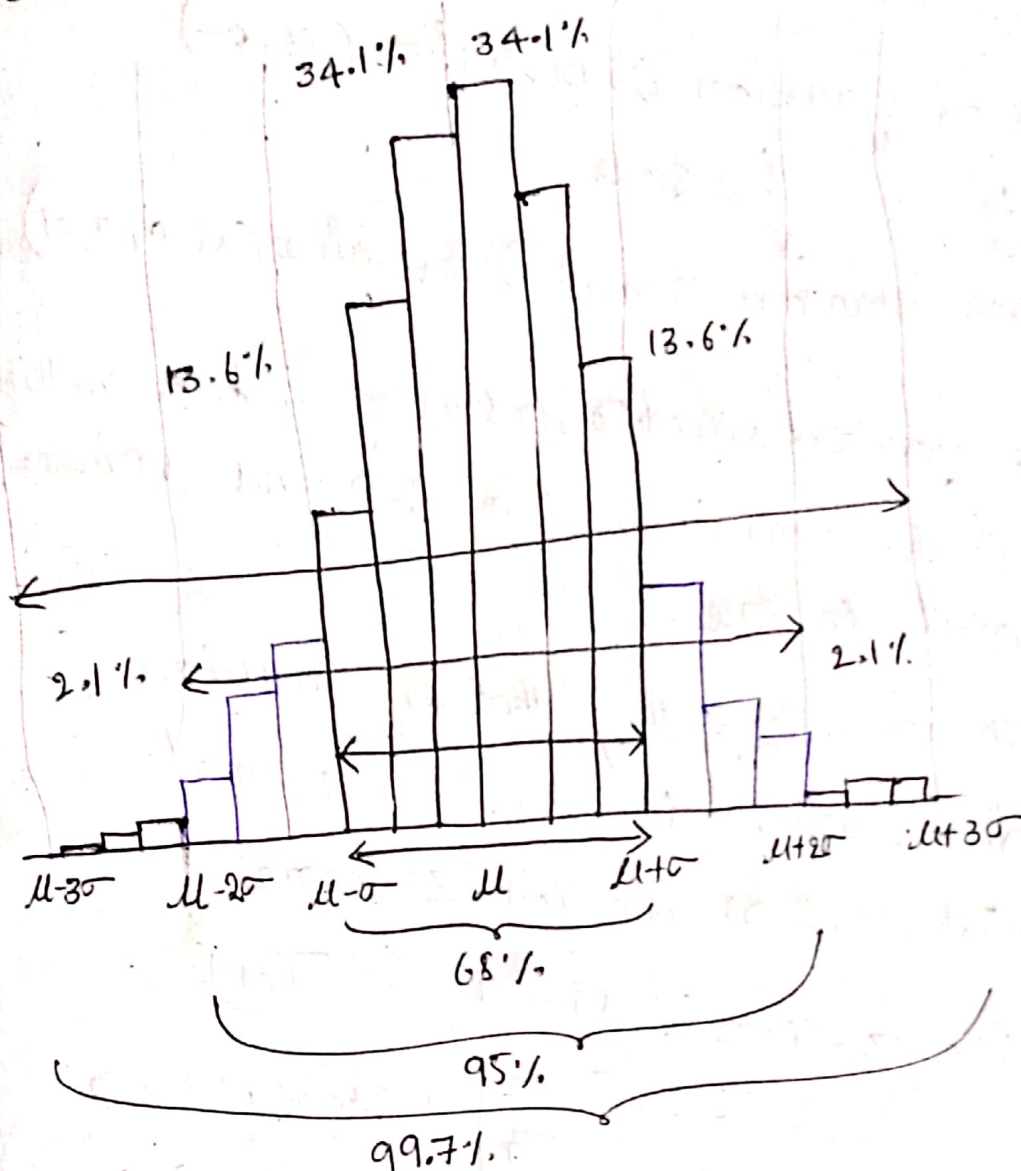
- \* Normal Distribution
- \* Standard Normal Distribution
- \* Z-Score
- \* Standardization and Normalization
- \* log Normal Distribution

### Gaussian / Normal Distribution!



- The Gaussian distribution is in the bell shaped curve.
- It is symmetrical
- The area covered under the curve is 1, that means 100%.
- why Gaussian (Normal distribution important)?
  - ⇓
  - Assumptions of data

# Empirical Rule of Normal distribution.



## Empirical rule:

68-95-99.71.

Q-Q plot  $\Rightarrow$  Distribution is Gaussian or not?

- \* The ~~in~~ the Empirical rule
- \* The sum of one standard deviation to right and one to left is covered 68% of data
- \* The sum of two standard deviation to right and two standard deviation is covered 95% of data.
- \* The sum of three standard deviation to right and three standard deviation is covered 99.7% of data.

## Standard Normal Distribution:

$x \approx$  Gaussian Distribution ( $\mu, \sigma$ )

$\Downarrow$

$\Downarrow$  Z-score

$y \approx$  Standard Normal Distribution ( $\mu=0, \sigma=1$ )

The standard normal distribution is here nothing but mean is equal to zero and variance equal to one.

→ to convert  $x$  to  $y$  that is Gaussian distribution to standard normal distribution we use Z-score

$$\text{i.e., Z-score} = \frac{x_i - \mu}{\frac{\sigma}{\sqrt{n}}} \quad \boxed{n=1}$$

$\frac{\sigma}{\sqrt{n}} \rightarrow \text{Standard error}$

$$\text{Z-score} = \boxed{\frac{x_i - \mu}{\sigma}} \quad (\because n=1)$$

Ex:  $x = \{1, 2, 3, 4, 5\} \rightarrow$  Gaussian distribution

$\Downarrow$

$$\mu = 3 \quad \sigma = 1.414$$

convert to SND

$$= \frac{1-3}{1.414} = -1.414$$

$$= \frac{2-3}{1.414} = -0.707$$

$$= \frac{4-3}{1.414} = \frac{1}{1.414} \approx 0$$

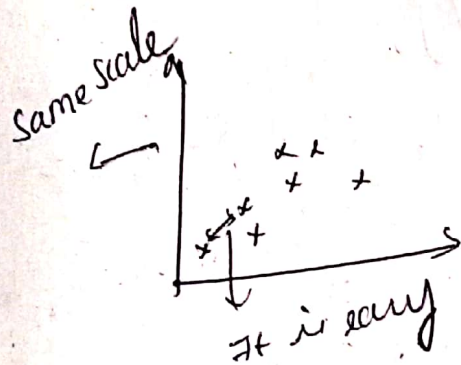
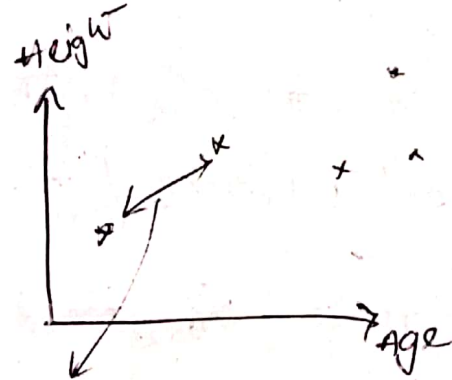
... so on ...

$$y = \{-1.414, -0.707, 0, 0.707, 1.414\}$$



why?

Age	weight	Height
24	72	150
26	78	160
32	84	165
33	92	170
34	87	150
28	83	180
29	80	175



## Feature Scaling

### 1. Standardization

$x \rightarrow$  Normal Distribution ( $\mu, \sigma$ )  
 $\Downarrow$  z-score

$y \rightarrow$  SND ( $\mu=0, \sigma=1$ )

why do we do this?

$\rightarrow$  To bring the features in the same scale.

## Normalization [0-1]

In standardization we have  $\mu=0$  and  $\sigma=1$  those are pre-fixed.

But in normalization we take the mean and variance [0-1]

Became to scale down

→ Min-max scalar [0-1]

$$x_{\text{scaled}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

$$x = \{1, 2, 3, 4, 5\}$$

$$x \Rightarrow y$$

$$1 \quad 0$$

$$2 \quad 0.25$$

$$3 \quad 0.5$$

$$4 \quad 0.75$$

$$5 \quad 1$$

$$\frac{4-1}{5-1} = \frac{3}{4}$$

$$\frac{3-1}{5-1} = \frac{2}{4}$$

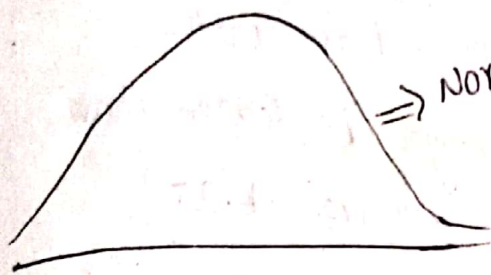
$$\frac{2-1}{5-1} = \frac{1}{4}$$

$$\frac{1-1}{5-1} = 0$$

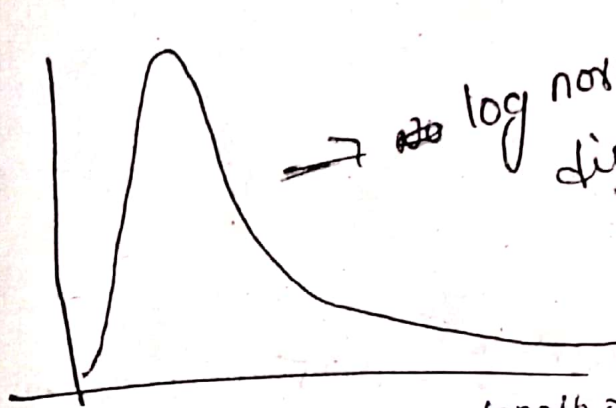
\* standardization is used in machine learning algorithms.

\* normalization is used in deep learning.

Log Normal Distribution



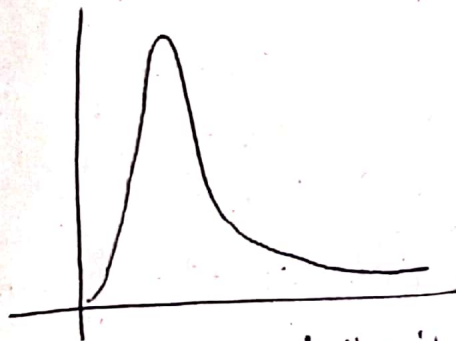
$\Rightarrow$  normal / Gaussian Distribution



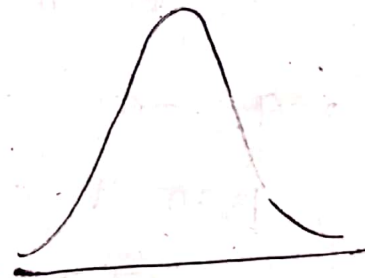
$\rightarrow$  ~~no~~ log normal distribution

$\rightarrow$  length of commands

\* To convert log normal Distribution to Gaussian

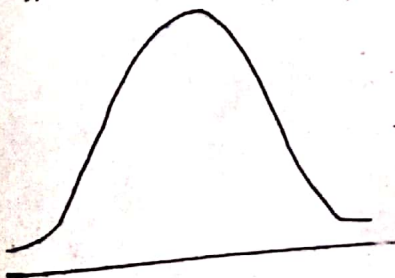


$$\Rightarrow y = \ln(x)$$

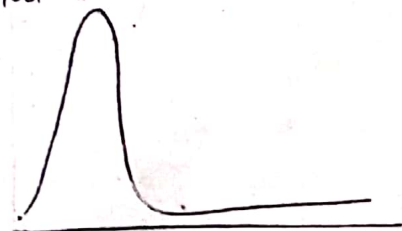


$x \approx$  log normal distribution

\* To convert Gaussian to log normal distribution

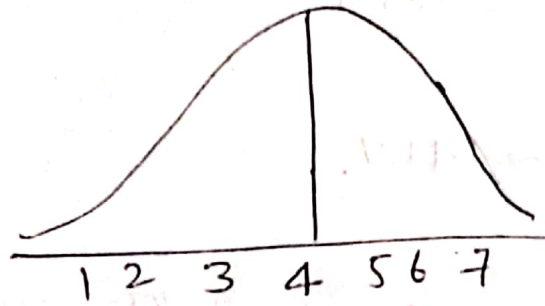


$$\Rightarrow x = \exp(y)$$



\*  $x = \{1, 2, 3, 4, 5\}$   $\mu = 4$ ,  $\sigma = 1$

$\Rightarrow$  what is the percentage of score that falls above 4.25?



$$z\text{-score} = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

z-table (Area under the curve)

1) See in google for z-table positive and negative values.

According to that values we can find the percentage of area covered under the curve.

$$\text{for } 0.25 \Rightarrow 0.598 \Rightarrow 0.598 \times 100\% = 59.8\%$$

$\Rightarrow$  what is the percentage of score that fall below 3.75

$$z\text{-score} = \frac{x_i - \mu}{\sigma} = \frac{3.75 - 4}{1} = -0.25$$

$$\Rightarrow -0.25 \Rightarrow 0.40 \Rightarrow 0.40 \times 100\% = 40\%$$

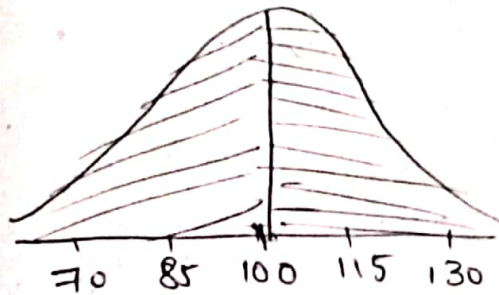


\* In India the average IQ is 100 with a standard deviation of 15. what is the percentage of population would you expect to have an IQ

(i) lower than 85

(ii) higher than 85

(iii) Between 85 and 100



(i) lower than 85

$$\begin{aligned} \mu &= 100 \\ \sigma &= 15 \\ \therefore Z\text{-score} &= \frac{x_i - \mu}{\sigma} \\ &= \frac{85 - 100}{15} = -1 \end{aligned}$$

$-1 \Rightarrow$  Z-negative score is 0.15  $\Rightarrow 0.15 \times 100\% = 15\%$

(ii) higher than 85

Z-score = +1  $\Rightarrow$  Z-positive score for +1 is 0.84  
 $\Rightarrow 0.84 \times 100\% = 84\%$

(iii) Between 85 and 100

first find lower  
 Subtract higher from higher after finding  
 Z-scores  $\Rightarrow 84\%$  covered higher than 85  
 $\Rightarrow 50\%$  covered lower than 100  
 $84\% - 50\% \Rightarrow 34\%$  between 85 & 100