

Lecture - 03

Distribution,

→ Normal (or) Gaussian distribution.

→ Standard Normal distribution.

→ Z score.

→ Log Normal distribution.

→ Bernoulli's distribution.

→ Binomial distribution.

Practical,

→ Mean, Median, Mode

→ Variance, Standard deviation.

→ Histogram, Pdf, Barplot, Violin plot.

→ IQR

→ Log Normal distribution.

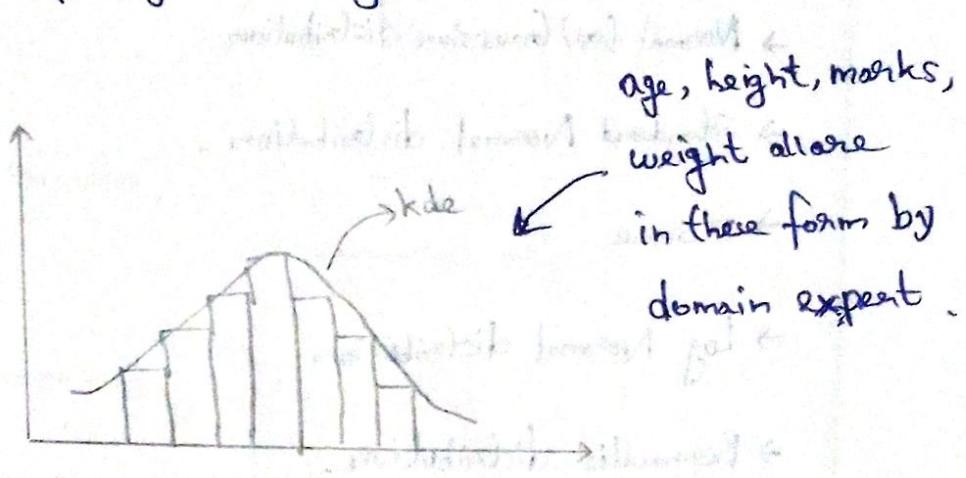
What is distribution?

How the data points are actually distributed
in the Visualized graph.

Example,

$$\text{age} = \{24, 26, 28, 15, 16, 17, 34, \dots\}$$

Plot histogram for age something look like,



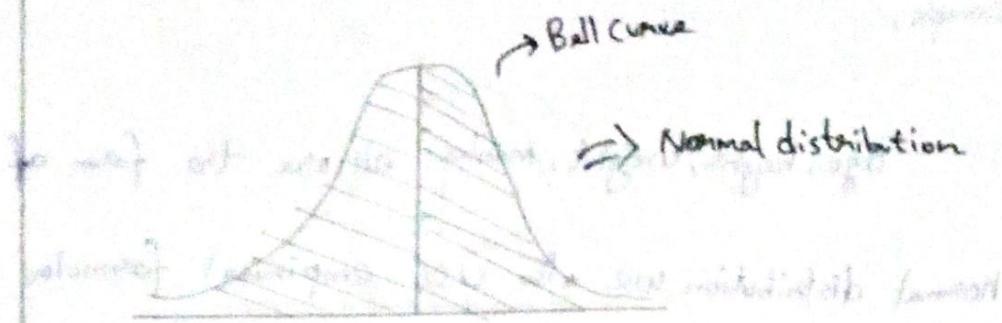
These type of distribution is known as Normal

(oo) Gaussian distribution.

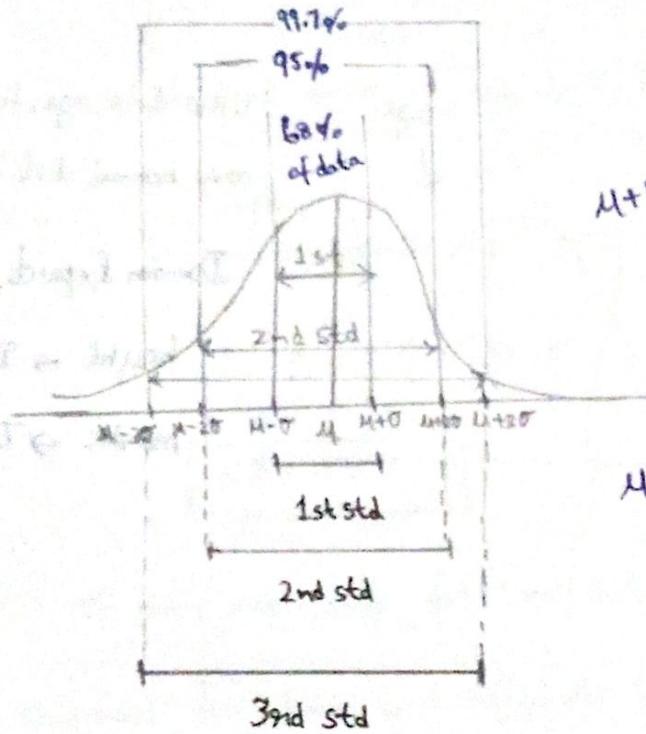
In world, Most of the data are follows the Gaussian distribution.

Gaussian / Normal Distribution:

→ It is the ~~histogram~~ distribution it looks like an bell - curve.



\hookrightarrow It may be Mean, Median (or) Mode



$\mu + \sigma \Rightarrow$ one standard deviation from right from the mean

$\mu - \sigma \Rightarrow$ one standard deviation left from the mean

Empirical formula : [This rule is applicable for normal dist]

This formula says the 68 - 95 - 99.7 % Rule

This rule is applied on the normal distribution.

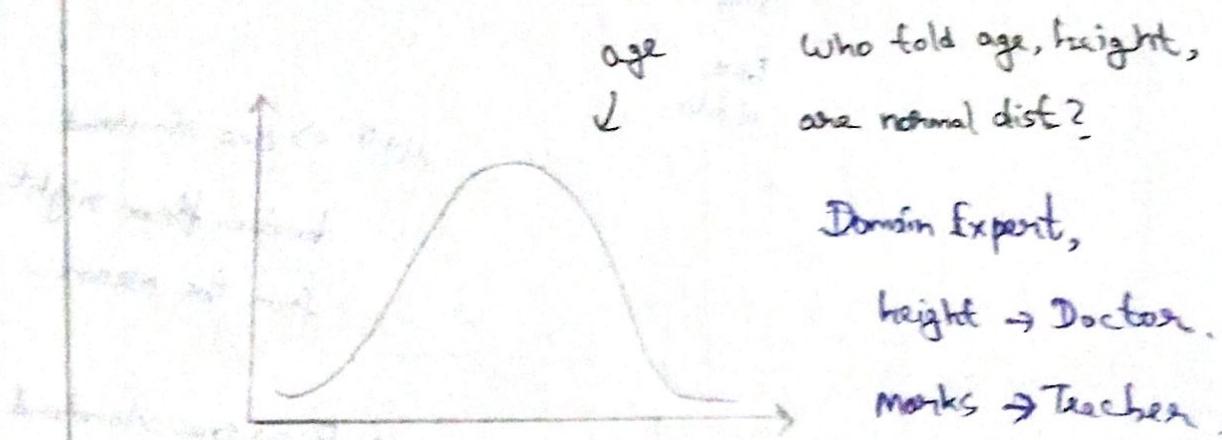
68% \rightarrow 1st std

~~2nd~~
95% \rightarrow 2nd std

99.7% \rightarrow 3rd std

Example,

age, height, weight, marks all are the form of normal distribution. we also use empirical formula for feature.

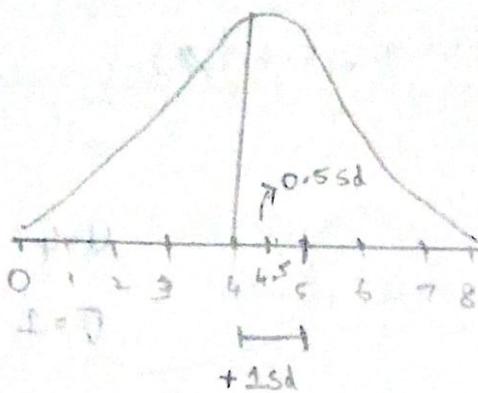


Z-Score:

\rightarrow Z-score refers to ~~how~~ find out the how much standard deviation of the data (σ) value from the mean.

\rightarrow Example,

$$\mu = 4, \sigma = 1$$



How much Standard deviation of 4.5 ?

Simply,

$$Sd = \sigma = 1$$

$$1 - 0.5 = 0.5$$

$$\boxed{Sd = +0.5}$$

This method

If only applicable for Small data points and Small no of standard deviation and difficult to find.

So, Using Z-Score to find out the Sd.

formula \Rightarrow

$$\boxed{Z\text{Score} = \frac{x_i - \bar{x}}{\sigma}}$$

To find 4.75,

$$Z\text{score} = \frac{4.75 - 4}{1} = +0.75 \text{ sd}$$

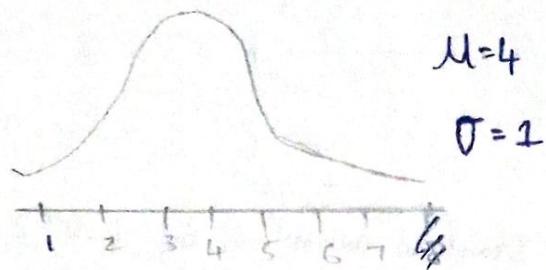
Suppose (-) to consider

the left sd from the Mean.

[obviously to find left or right sd use z-score]

Example,

$$\text{data} = \{1, 2, 3, 4, 5, 6, 7, 8\}$$



To find Out Z-score for every element,

$$Z(1) = \frac{1-4}{1} = -3$$

$$Z(4) = \frac{4-4}{1} = 0$$

$$Z(2) = \frac{2-4}{1} = -2$$

$$Z(5) = \frac{5-4}{1} = 1$$

$$Z(3) = \frac{3-4}{1} = -1$$

$$Z(6) = \frac{6-4}{1} = 2$$

$$Z(7) = \frac{7-4}{1} = 3$$

$$\text{data} = \{1, 2, 3, 4, 5, 6, 7\}$$

↓ After Applying Z-score for each.

$$\text{data} = \{-3, -2, -1, 0, 1, 2, 3\}$$

-3 sd -2sd +3 sd
from from Mean
Mean

+3 sd
from
Mean

Note that,

$$\text{data} = \{-3, -2, -1, \underline{0}, 1, 2, 3\}$$



$$\text{definitely } \mu=0, \sigma=1$$

If $\mu=0, \sigma=1$, then the datapoints are fall in Standard normal distribution.

The given data is the distribution of Normal distribution, then we apply Z-score. The data are converted to the Standard Normal distribution.

Standard Normal distribution:

Properties

$$\mu=0, \sigma=1$$

$$\text{Random Variable } (X) = X \sim \text{SND}(\mu=0, \sigma=1)$$

Real time Example:

Consider the dataset with feature age, salary, weight.

different units

<u>Age</u>	(years)	<u>Salary</u>	(Rs)	<u>Weight</u>	(kg)
24		40k		70	
25		80k		80	
26		60k		55	
27		170k		45	

APPLY Z-score when,

$\mu = 0, \sigma = 1$

This applying z-score process is known as

Standardization.

Standardization: [Uses z-score]

Standardization is an important technique that is mostly performed pre-processing step. To standardize the range of features of an input data set.

Many of algorithm need Standardization because of # features have different units to scale.

Example,

Weight

Height

Consider Weight and height in the dataset.

To build machine learning model, Height is dominated with weight. Maybe height will be broad range of value. To avoid So, we need to standardize the value for height and weight. Range: (-) to (+)

Normalization: [uses Min Max Scaler]

Normalization is also method for used in feature Scaling and pre-processing steps. It Using ~~for~~ Min-Max Scaler. [There are so many formula applied in Minmax Scaler based

Range = 0 to 1 on needs]

Example, In CNN (Image Classification)



Every pixel range from 0 to 255

Using minmax scalar to convert range to 0 to 1.

(00)

$$\frac{0}{255} = 0, \frac{1}{255} = \dots, \frac{255}{255} = 1 \quad [0, \dots, 1]$$

Practical Examination:

ODI Series for Cricket of 2020 and 2021.

2020

Series average score = 250

Standard deviation = 10

Team final score = 240

2021

Series average score = 260

Standard deviation = 12

Team final score = 245

Compare to both the scores in which year team final score was better (or) perform well?

Ans,

In 2020,

$$Z\text{ score} = \frac{240 - 250}{10} = \frac{-10}{10} = -1$$

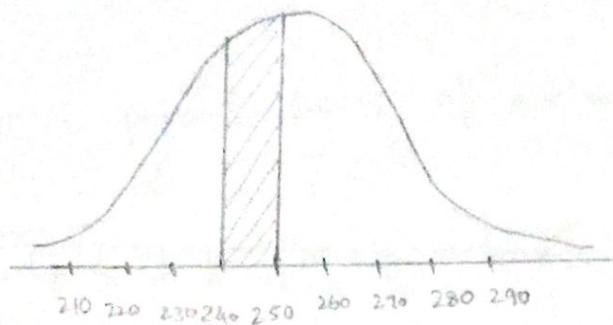
Z score of 2020 = -1

In 2021,

$$Z\text{ score} = \frac{245 - 260}{12} = \frac{-15}{12} = \frac{5}{4} = -1.25$$

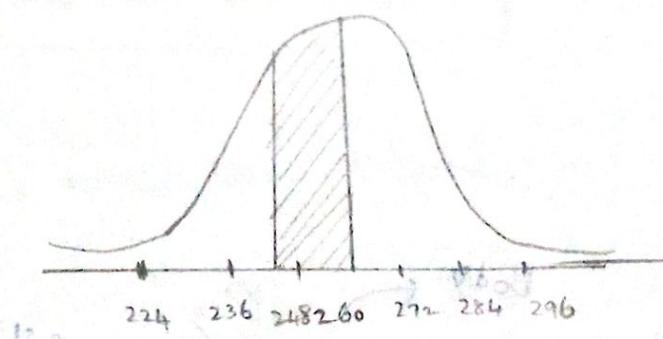
In 2020,

$$\mu = 250 \quad x_i = 240 \quad \sigma = 10$$



In 2021,

$$\mu = 260, \quad x_i = 245, \quad \sigma = 12$$



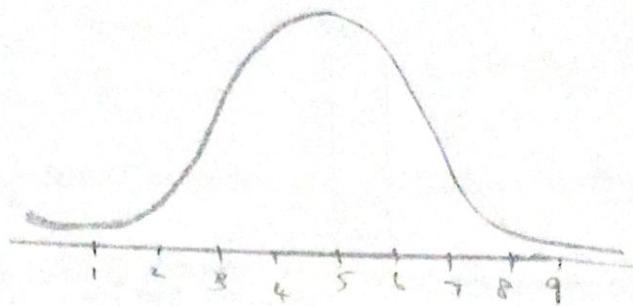
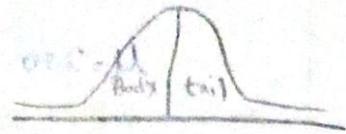
Based on the area spread in 2021 is high compared to 2020. It means the data ranges may be high between the area. So, definitely 2021 was better score compare than 2020.

Answer is correct

Prakash

Stats Interview Question:

Q1 = 3 Q3 = 7

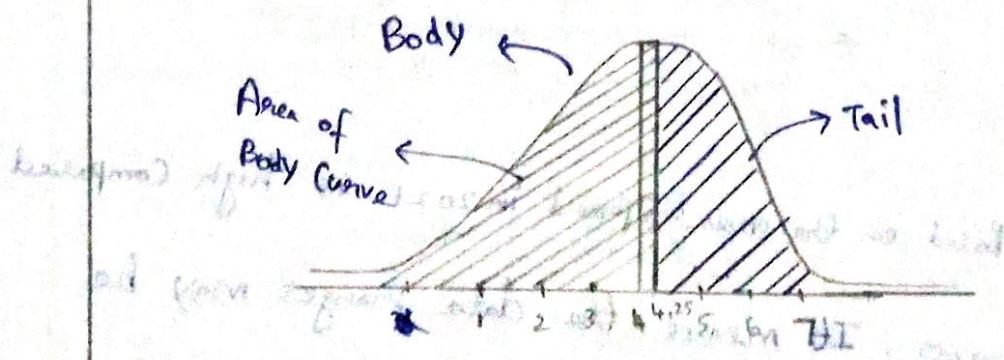


①

Question:

What Percentage of Scores fall above 4.25?

Consider $\mu = 4, \sigma = 1$



To find Scores fall above 4.25?

• Once with integral areas
Two ways,

①. First using basic Simple

Percentage,

~~The Z-Score has two different methods. It depends very less on informed result. Only support from some data.~~

$$\text{Percentage} = \frac{4.25}{7} \times 100 = 60.71428\% \quad \text{data-} \\ = 60\%$$

4.25 is present 60% of entire data.

Our Question is to fall above 4.25?

So Simple,

$$100\% - 60\% = 40\%$$

Scores above 4.25 = 40%

② Second Method Using Z-score,

$$Z\text{score} = \frac{4.25 - 4}{1} = 0.25$$

See Zscore table with respect to ~~these~~.

$$0.2 \text{ and } 0.05 = 0.59871 \quad \text{left area region}$$

We need above 4.25, ~~area~~ - 1

$$1 - \text{left area} = 1 - 0.59871$$

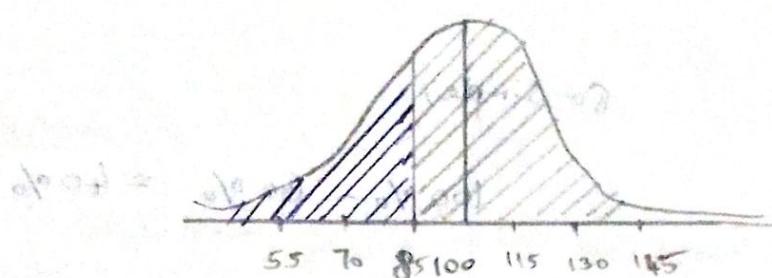
$$= 0.40129 \Rightarrow 40\%$$

(2)

Question:

$$\sigma_{\text{population}} = \text{SD} \times \frac{2\pi}{3} = \text{approx}$$

In India the average IQ is 100, with a standard deviation of 15. What percentage of the population would you expect to have an IQ lower than 85?



$$\mu = 100, \sigma = 15$$

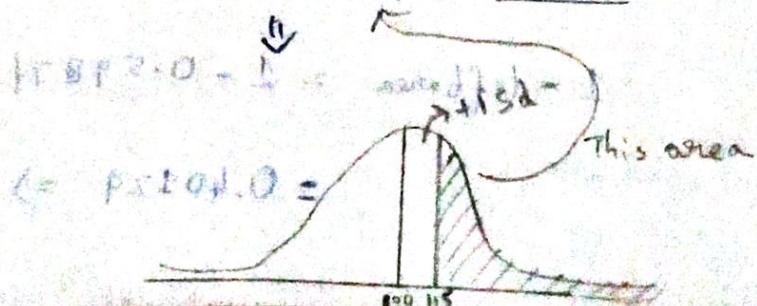
$$Z\text{ score} = \frac{85 - 100}{15} = \frac{-15}{15} = -1$$

In Z-table (-1.0) is 0.15866 [direct look in
(or) 0.84134 1 - 0.84134 0.15866 Ztable]

See (1.0) in Z-table is 0.84134

So,

$$1 - 0.84134 + \underline{0.15866}$$



In Gaussian distribution, Both the sides are equal.
The same way both the SD size area also same.

So, The area is 0.15866 is nothing but 15.8%.

Real time practical Using Python,

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
%matplotlib inline
```

```
import Seaborn as sns
```

```
import Statistics
```

```
df = sns.load_dataset('tips')
```

```
df.head()
```

```
# Mean, mode, median
```

```
np.mean(df['total_bill'])
```

```
np.median(df['total_bill'])
```

```
Statistics.mode(df['total_bill']) # well suited for Categorical  
Value
```

```
# Visualize boxplot
```

```
Sns.boxplot(df['total_bill'])
```

Visualize histplot

Sns.histplot(df['total_bill'])

Sns.histplot(df['total_bill'], kde=True) # Pdf applied.

Pdf used to whether the distribution is gaussian
(or) not.

Sns.cookplot(df[parameter])

Find Percentile of feature values.

$$IQR = Q_3 - Q_1$$

np.percentile(df['total_bill'], [25, 75])

np.percentile(df['total_bill'], [25, 99])