

Agenda

- ① Percentiles and Quartiles ✓
- ② 5 Number Summary {Outliers} ✓
- ③ Box plot ✓
- ④ Covariance And Correlation }.
- ⑤ Probability distribution function
- ⑥ Different types of distribution

① Percentiles And Quartiles [GATE, CAT]

Percentage : 1, 2, 3, 4, 5, 6

$$\% \text{ of numbers that are odd} = \frac{3}{6} = \frac{\text{No. of odd numbers}}{\text{No. of total no.}}$$

$$= \frac{1}{2} = 50\%$$

Percentiles :

Defn : A percentile is a value below which a certain percentage of data points lie.

$n = 15$

$$X = \{2, 3, 3, 4, 6, 6, 6, 7, 8, 8, 9, 9, 10, 11, 12\}$$

$$\text{Percentile Rank of } \underline{\underline{10}} = \frac{\text{Value} \# \text{ of values below } 10}{n} * 100$$

$$= \frac{7 + 8}{15} \times 100 = 80 \text{ percentile}$$

80 percentile = 80% of the distribution fall below the value of 10 //

Here are the steps on how to calculate the percentile of a value:

① What value exists at 25 percentile?

Collect the data set.

Arrange the data set in ascending order.

Determine the total number of observations.

Identify the data value for which you are interested to find the percentile.

Count the number of data values that are less than the above value.

Divide the number from Step 5 by the number from Step 3 to find the percentile of the given data value.

$$\text{Value} = \frac{\text{Percentile}}{100} * (n+1)$$

$$= \frac{18}{100} * 12 = \boxed{4} \text{ Element } = 4$$

20/41

$$X: \{ \begin{matrix} \downarrow & \downarrow & \downarrow & \downarrow \\ 2, 3, 3, 4, 6, 6, 6, 7, 8, 8, 9, 9, \boxed{10}, 11, 12 \end{matrix} \}$$

$$\frac{4+1}{2} = \underline{\underline{5}}$$

4.5

If value will come in decimal, then take avg of that no. before decimal and the next no.

then that's no. will be the value

Quartiles

① Q1 → 25 percentile

Q2 → Median → 50 percentile

Q3 → 75 percentile

② 5 Number Summary

(1) Minimum

(2) First Quartile (25 percentile) (Q1)

(3) Median (Q2)

(4) Third Quartile (75 percentile) (Q3)

(5) Maximum

After removing the outlier

Meaning of taking this lower fence and higher fence is that below lower fence everything will be outlier and above higher fence everything will be outliers.

These fence only use for checking outliers.

$$X = \{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \boxed{12, 9} \}$$

[Lower Fence \longleftrightarrow Higher Fence]

$$\text{Lower Fence} = Q_1 - 1.5(IQR)$$

$$\text{Inter Quartile Range} = Q_3 - Q_1$$

$$\text{Higher Fence} = Q_3 + 1.5(IQR)$$

$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\boxed{12, 9}}\}$. ↓ Outlier

$$Q_1 = 25^{\text{percentile}} = \frac{25}{100} * 20 = 5^{\text{th value}} = 3$$

$$Q_3 = 75^{\text{percentile}} = \frac{75}{100} * 20 = 15^{\text{th value}} = 7$$

$$IQR = 7 - 3 = 4$$

$$\text{Lower Fence} = Q_1 - 1.5(IQR)$$

$$= 3 - 1.5(4)$$

$$= 3 - 6$$

$$= -3$$

$$\text{Higher Fence} = Q_3 + 1.5(IQR)$$

$$= 7 + 1.5(4)$$

$$= 7 + 6$$

$$= 13$$

$$[-3, 13].$$

$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\boxed{12, 9}}\}$. ↑ Outlier

Box plot [To visualize Outliers]

5 number summary is used for Box plot

Box plot

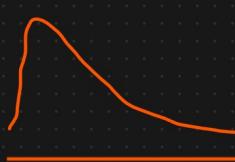
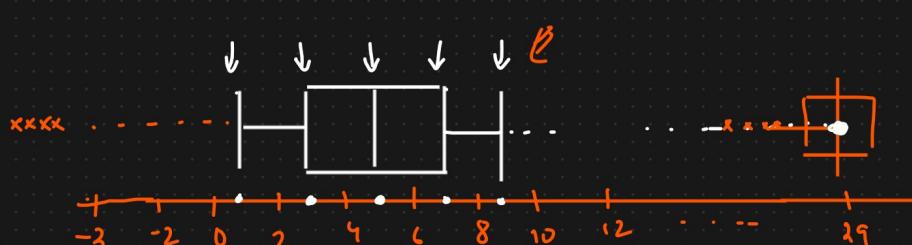
$$\textcircled{1} \text{ Minimum value} = 1$$

$$\textcircled{2} \quad Q_1 = 3$$

$$\textcircled{3} \quad \text{Median } Q_2 = 5$$

$$\textcircled{4} \quad Q_3 = 7$$

$$\textcircled{5} \quad \text{Maximum} = 9$$

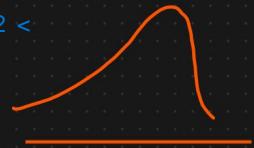


$$Q_3 - Q_2 > Q_2 - Q_1$$

mean > median > mode

Q2 is left side as more data is towards right side. So $Q_3 - Q_2$ should be higher than $Q_2 - Q_1$.

Here $Q_3 - Q_2 < Q_2 - Q_1$



Internal Assignment

$$-5+3 = 1 \quad \uparrow \overline{2}$$

$$y = \{-13, -12, -6, \boxed{5}, 3, 4, 5, 6, 7, 7, 8, \boxed{10}, \boxed{10}, 11, 24, 55\}.$$

[lower fence \longleftrightarrow higher fence]

$$Q_1 = \frac{21}{4} \times 17^4 = 4.25$$

$$\text{lower fence} = -1 - 1.5(10 + 1)$$

$$Q_3 = \frac{75}{100} \times 17 = 12.75$$

$$= -1 - 1.5(11)$$

$$= -17.5$$

$$[-17.5, 26.5].$$

$$\text{higher fence} = 10 + 1.5(10 + 1)$$

↓
55 is an outlier

$$= 10 + 16.5$$

$$= 26.5$$

\equiv

$$Z = \{1, 2, 4, 6, 7, 12, 18, 34, \boxed{77}, \boxed{66}, 108, 99, 14\}.$$

$$Q_1 = 5 \\ \equiv$$

$$Q_3 = 71.5 \\ \equiv$$

$$[-94.75, 171.25] \\ \equiv$$

$$\{1, 2, 4, 6, 7, 12, 18, 34, \boxed{66}, \boxed{77}, 99, 108\} \\ 14, \\ 153$$

$$Q_3 = \frac{75}{100} \times 14 = \frac{42}{4} 2 + 10.5$$

$$Q_3 = \frac{66 + 77}{2} = 71.5 \\ \equiv$$

Covariance And Correlation

[Relationship between X and Y]

X	Y
2	3
4	5
6	7

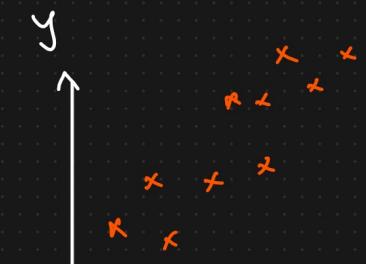
→	X↑	Y↑
→	X↓	Y↑
→	X↑	Y↓

IS
Household

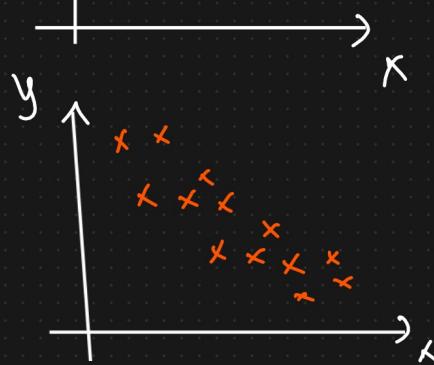
Size of
house

Price

$$8 \quad 9 \quad \rightarrow \boxed{x \downarrow \quad y \downarrow}$$



$$\boxed{\begin{matrix} x \uparrow & y \uparrow \\ x \downarrow & y \downarrow \end{matrix}}$$



$$\boxed{\begin{matrix} x \uparrow & y \downarrow \\ x \downarrow & y \uparrow \end{matrix}}$$

① Covariance

$$\text{Cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1} \quad \text{Var}(x) = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

After applying above formula if ---

$$= \sum_{i=1}^n \frac{(x_i - \bar{x}) \times (x_i - \bar{x})}{n-1}$$

$$\text{Cov}(x, y) \rightarrow \boxed{\begin{matrix} x \uparrow y \uparrow \\ x \downarrow y \downarrow \end{matrix}} \Rightarrow \text{tve Covariance}$$

$$\boxed{\text{Var}(x) \subset \text{Cov}(x, x)}$$

$$\boxed{\begin{matrix} x \uparrow y \downarrow \\ x \downarrow y \uparrow \end{matrix}} \Rightarrow \text{-ve Covariance.}$$

Variance of x is covariance of x, x means relationship of x with x itself is covariance of (x, x)

$$\text{Cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$= \left[(2-4)(3-5) + (4-4)(5-5) + (6-4) \times (7-5) \right]$$

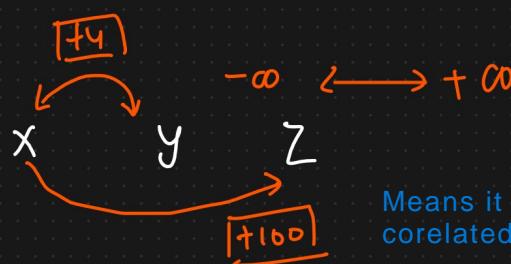
$$\begin{array}{cc} x & y \\ \rightarrow 2 & 3 \\ \rightarrow 4 & 5 \\ \underline{6} & \underline{7} \\ \bar{x} = 4 & \bar{y} = 5 \end{array}$$

$$= \frac{4+0+4}{2} = \frac{8}{2} = 4 \quad \text{tve Covariance}$$

X and Y are having a positive covariance

Advantages

- ① Relationship between X & Y



Disadvantage

- ① Covariance does not have a specific limit value

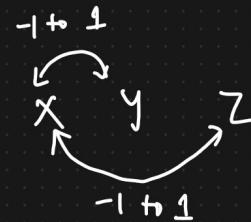
It can be any +ve value or any -ve value but can't be any limit or strength that is why there is no limit in covariance

Means it is not saying if X is highly correlated with Y or Z.

- ② Pearson Correlation Coefficient Range $[-1 \text{ to } 1]$

For linear dataset

$$r_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y}$$



- ③ The more the value towards +1 the more the correlated it is
 ④ The more the " " -1 " " " -ve "

X Y 0.6

So from here we can say that X is highly correlated with Z than Y with the help of Pearson correlation function.

X Z 0.7

- ⑤ Spearman Rank Correlation

For non linear dataset will also be captured

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sqrt{R(x) \cdot R(y)}}$$

X	Y	R(x)	R(y)
5	6	3	1
7	4	2	2
8	3	1	3
1	1	5	5
2	2	4	4

Created rank wise table which are bigger no.

Feature Selection

How these value correlated with price

+ve	+ve	+ve	≈ 0	-ve	$\frac{0/p}{\text{Price}} \uparrow$
Size of house	No. of rooms	location	No. of people staying	height	

Probability Distribution Function And Probability Density Function

Probability Mass function

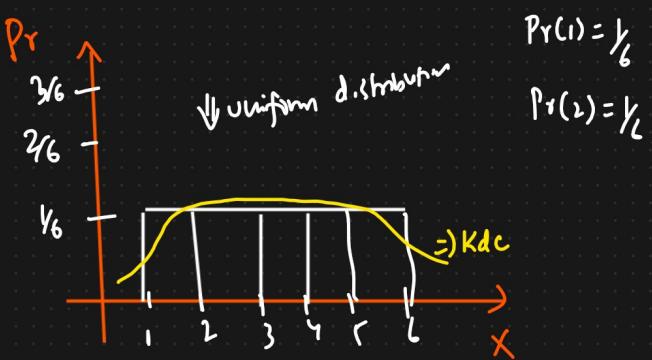
Probability Distribution Function

- ① Probability density fn
- ② Probability mass fn ✓
- ③ Cumulative distributn fn.

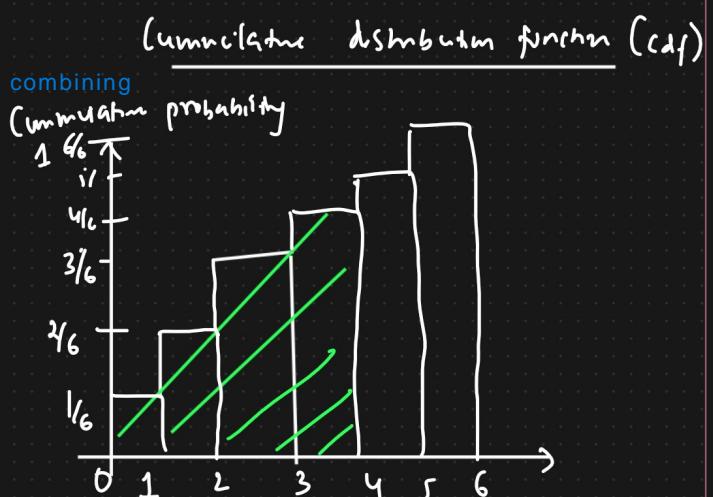
① PMF ↴

① Discrete Random Variable ⇩

Eg: Rolling a dice {1, 2, 3, 4, 5, 6}



$$\Pr(1 \text{ or } 2) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

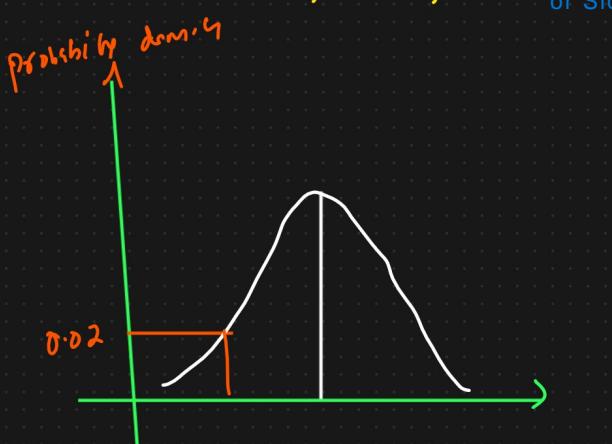
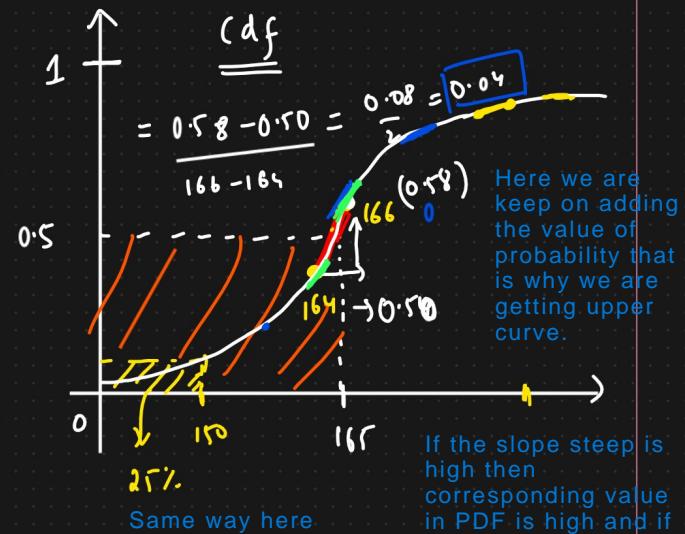
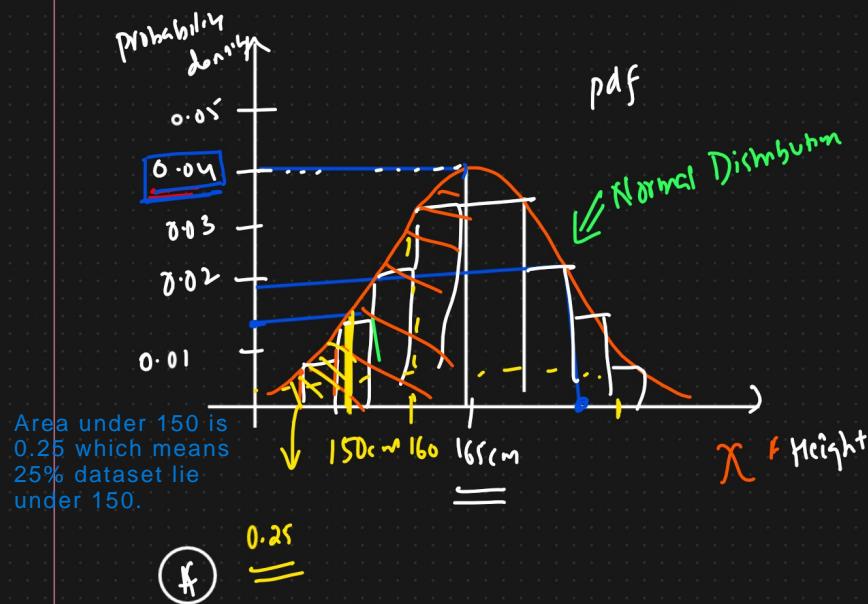


$$\Pr(X \leq 4) = \Pr(X=1) + \Pr(X=2) + \Pr(X=3) + \Pr(X=4)$$

= 0/p

② Probability Density Function

① Distribution of continuous Random Variable



Different types of Distribution

① Normal / Gaussian Distribution \rightarrow pdf

② Standard Normal Distribution \rightarrow pdf

③ Log Normal Distribution \rightarrow pdf

④ Power law Distribution \rightarrow pdf

⑤ Bernoulli Distribution \rightarrow pmf

⑥ Binomial Distribution \rightarrow pmf

⑦ Poisson Distribution \rightarrow pmf

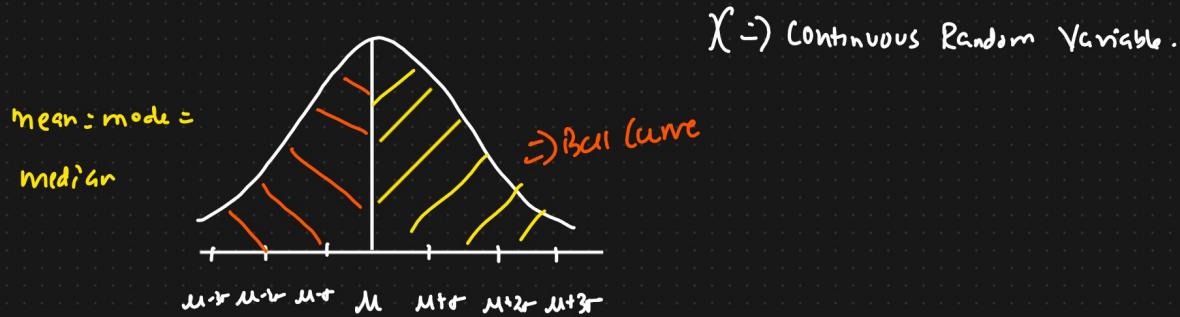
⑧ Uniform Distribution \rightarrow Discrete \rightarrow pmf
 \rightarrow Continuous \rightarrow pdf

⑨ Exponential Distribution. \rightarrow pdf

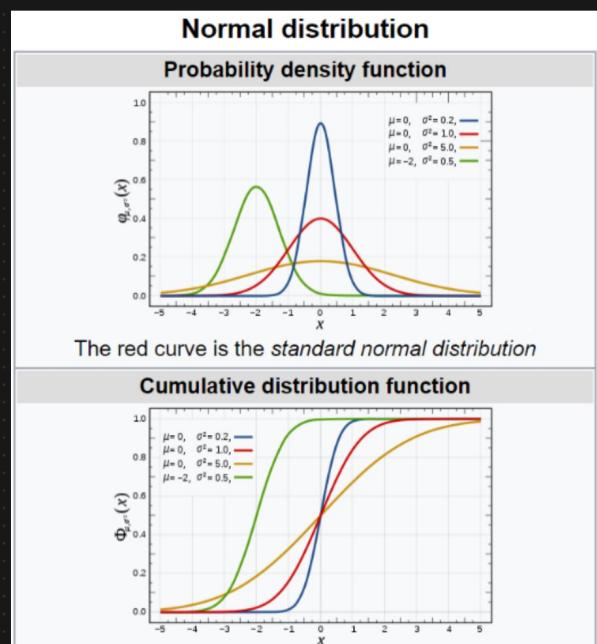
⑩ CHI SQUARE Distribution \rightarrow pdf

⑪ F Distribution \rightarrow pdf.

⑫ Normal/Gaussian Distribution



Eg:- Height, weight, age, IRIS dataset



$$X \sim N(\mu, \sigma^2)$$

Support parameters $\mu = \text{mean}$
 $\sigma^2 = \text{variance}$

$$\text{PDF} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

100 datapoint

Empirical Rule

68 - 95 - 99.7% Rule

$$X = \{ \quad \}$$

mean = mode =

median

The empirical rule can be used to make inferences about the probability of a particular value occurring in a normally distributed population. For example, if we know that the average height of a population of adults is 170 centimeters and the standard deviation is 5 centimeters, we can use the empirical rule to say that 68% of the adults in the population will have heights between 165 and 175 centimeters.



$$m-3\sigma \quad m-2\sigma \quad m-\sigma \quad m \quad m+\sigma \quad m+2\sigma \quad m+3\sigma$$

First Standard deviation range -- 68% data lie

68%

Second SD range -- 95% data lie

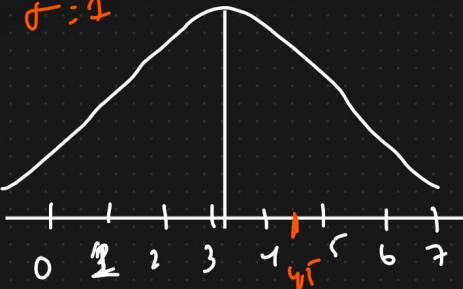
Third SD range -- 99.7 % data lie

Standard Normal Distribution

X

$$\mu = 3 \quad \sigma = 1$$

$$\mu = 3 \quad \sigma = 1$$



↳ Standard Normal Distribution

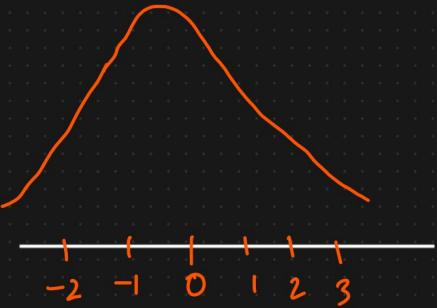
Then we get

Applied some

Transformation

$$\Rightarrow$$

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



$$\downarrow \\ Z\text{-score} = \frac{x_i - \mu}{\sigma}$$

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

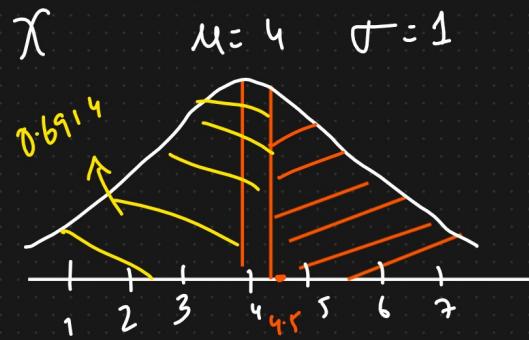
$$= \frac{2 - 3}{1} = -1$$

Z-score tells you about a value how many Standard deviation away from the mean

$$\frac{3-3}{1} = 0$$

$$\boxed{4.5} \Rightarrow \frac{4.5 - 3}{1} = 1.5$$

In other words, if we want to find for any value that it is how many Standard deviation Away from the mean. Then Z-score is used.



What is the percentage of scores above 4.5?

6.52

Simply mean what is the area above 4.5 in the graph

\Rightarrow Z-table

$$Z\text{-score} = \frac{4.5 - 4}{1} : 0.5 \Leftarrow \text{Means 0.5 SD Away from the mean}$$

Area under curve = $1 - 0.6914$
Check from z table from google PDF and check for +ve SD
0.5 and it will give area under curve left side that is why doing minus for area above 4.5 score here.

$$= 0.3086 // = 30.86\%$$