

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

Studying these just to determine and visualize outliers in the data set.

All these work in sorted data

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.

Eg: The CAT cutoff for IIM is typically from 99 to 100 percentile. So only top 1 percentage of population is able to make into IIMs.

$$n = 15$$

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

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

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

$$75.8$$

What 80 percentile means

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

② What value exists at 25 percentile?

Formula

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

$$= \frac{18}{100} * 164 = \boxed{4} \text{ Element } = 4$$

Fourth element is 4 after arranging number in ascending order.

$$\begin{aligned} X: & \{2, 3, 3, 4, 6, 6, 6, 7, 8, 8, 9, 9, \boxed{10}, 11, 12\} \\ & \begin{array}{ccccccc} \downarrow & \downarrow & \downarrow & \downarrow & & & \\ 1^{\text{st}} & 2^{\text{nd}} & 3^{\text{rd}} & 4^{\text{th}} & & & \\ \uparrow & \uparrow & \uparrow & \uparrow & & & \end{array} \\ & \frac{4+6}{2} = \underline{\underline{5}} \end{aligned}$$

If 25 percentile value had came in decimal EG: 4.5 then simply take average of 4th and 5th element (after arranging data in ascending order) and the result will be value at 25 percentile.

Why we are learning this?

Ans: Because later these will be used to calculate the outliers.

4.5

Quartiles

① Q1 → 25 percentile

Q2 → Median → 50 percentile

Q3 → 75 percentile

② 5 Number Summary

5 number summary is represented using Box Plot

(1) Minimum

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

(3) Median (Q2)

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

(5) Maximum

Remove The Outliers

Outlier is a value that is completely unique or not following a general trend.

$$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \boxed{29}\}$$

Outlier

[Lower Fence \longleftrightarrow Higher Fence]

Through observation we can tell in above 29 is outlier. But it's important to mathematically prove the same. Below we have tried to obtain outliers by drawing higher and lower fence, above or below which data points will be treated as outliers.

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

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

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

$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\boxed{29}}\}$. ↓ Outlier

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

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

$$IQR = 7 - 3 = 4$$

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

$$= 3 - 1.5(4)$$

$$= 3 - 6$$

$$= -3$$

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

$$= 7 + 1.5(4)$$

$$= 7 + 6$$

$$= 13$$

Note upper and lower fence value are inclusive and should not be considered outlier. Point above and below should be considered outliers.

$$[-3, 13]$$

Since 29 is above higher fence, it is considered as outlier.

There is no data point below lower fence.

$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\boxed{29}}\}$. ↑

Box plot [To visualize Outliers]

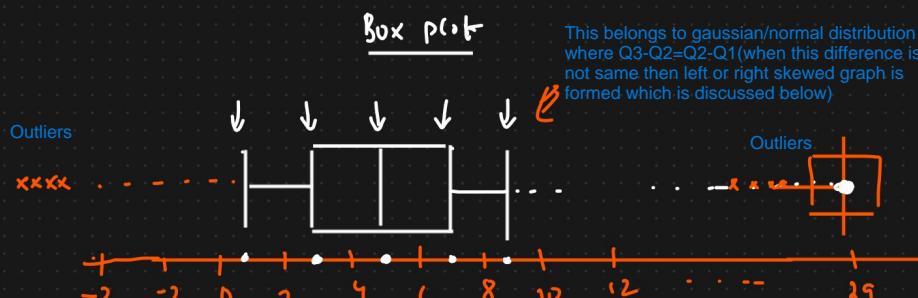
$$\textcircled{1} \text{ Minimum value} = 1$$

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

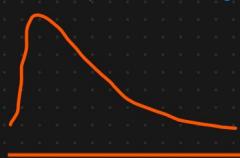
$$\textcircled{3} \text{ Median } Q_2 = 5$$

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

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



Interview question: Box plot for left and right skewed



mean > median > mode

Box plot

This belongs to gaussian/normal distribution where $Q_3 - Q_2 = Q_2 - Q_1$ (when this difference is not same then left or right skewed graph is formed which is discussed below)



$$Q_3 - Q_2 > Q_2 - Q_1$$

Case 1 : for right skewed graph box plot will be shifted left (ie, median = Q2 will shift left)

Means: $Q_3 - Q_2 > Q_2 - Q_1$

Case 2 : for left skewed graph box plot will be shifted right (ie, median = Q2 will shift right)

Means: $Q_3 - Q_2 < Q_2 - Q_1$

Internal Assignment

$$-5+3 = 1$$

$\uparrow \overline{2}$

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

Y is a random variable.

[lower fence \longleftrightarrow higher fence]

$$Q_1 = \frac{21}{100} \times 17^4 = 4.25$$

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

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

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

$$= -17.5$$

$$[-17.5, 26.5].$$

↓

55 is an outlier

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

$$= 10 + 16.5$$

$$= 26.5$$

=====.

Different question where Z is random variable

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

$$Q_1 = 5 \\ \underline{\underline{=}}$$

$$Q_3 = 71.5 \\ \underline{\underline{=}}$$

$$[-94.75, 171.25] \\ \underline{\underline{=}}$$

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

$$Q_3 = \frac{77}{100} \times 14 = \frac{42}{4} 2 + 10.5 \\ \underline{\underline{=}}$$

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

Covariance And Correlation

[Relationship between X and Y]

X	Y
2	3
4	5
6	7

→	X↑	Y↑
→	X↓	Y↑
→	X↑	Y↓

IS
handed

Size of
house

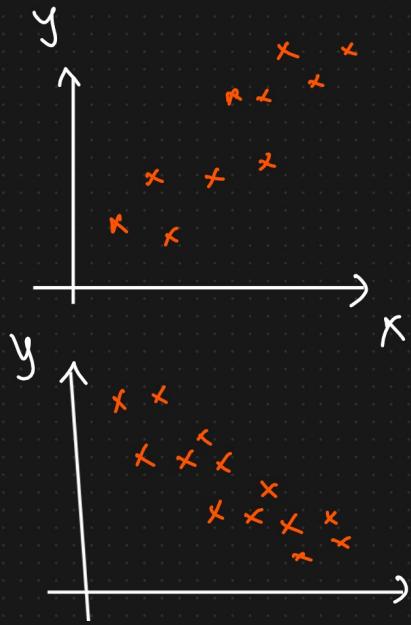
Price

All possible relationships between 2 variables. 1 more can be added that is by changing x = y)

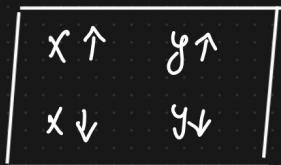
8

9

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



Case 1: Direct relationship
By increasing x, y increases and by decreasing x, y decreases.



Case 2: Indirect/inverse relationship
By increasing x, y decreases and by decreasing x, y increases.



Our aim is to get some values out of these relationships OR we just want to establish this whole thing mathematically. So, the first technique for doing this is Covariance.

① Covariance

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

How to remember. This formula is similar to the formula for variance just remove square(power) and add y - ybar

Covariance simply means relationship of one with another random variable. It can be +ve or -ve.

$$\text{Cov}(x, y)$$

Case 1:
For direct
relationship +ve Cov

$$\begin{array}{l} x \uparrow \quad y \uparrow \\ x \downarrow \quad y \downarrow \end{array}$$

\Rightarrow +ve
Covariance

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

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

Covariance of x
with itself = Cov(x,x)

\Downarrow

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

Case 2:
For indirect/inverse
relationship -ve Cov

$$\begin{array}{l} x \uparrow \quad y \downarrow \\ x \downarrow \quad y \uparrow \end{array}$$

\Rightarrow -ve
Covariance.

Var(X) means how well X is spread.

Cov(X,X) means relationship of X with itself

Below is numerical example of Cov(x,y)

X	Y
2	3
4	5
6	7
$\bar{x} = 4$	$\bar{y} = 5$

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

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

$$= \frac{2+0+4}{2} = \frac{8}{2} = 4 = +ve \text{ covariance}$$

X and Y are having a positive covariance

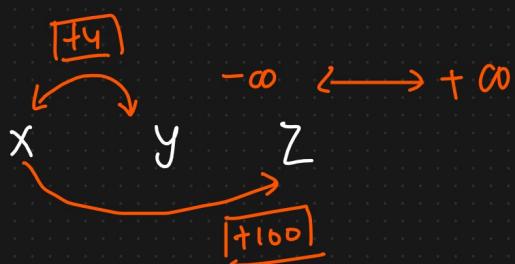
This is interview question:
What is relationship b/w variance and covariance.

Ans: Variance of any random variable X is equal to Covariance of that random variable X with itself.

i.e; $\text{Var}(X) = \text{Cov}(X, X)$

Advantages

① Relationship between x & y



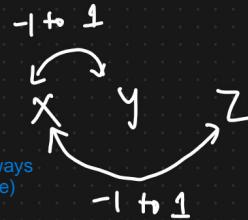
① Pearson Correlation Coefficient

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

Read as row of x and y

$[-1 \text{ to } 1]$

O/P of pearson corr coeff will always range from -1 to +1 (both inclusive)



Cov(x,y) = Covariance of x and y
 σ_x = standard deviation of x
 σ_y = standard deviation of y

- ④ The more the value towards +1 the more the correlated it is
- ② The more the " " " -1 " " " -ve " "

Pearson corr coeff of x and y = 0.6
 Pearson corr coeff of x and z = 0.7

Now, we can say that x is highly correlated with z as compared to correlation of x with y. This is something which we were not able to determine using Covariance alone.

↳ **Issue**

--> Pearson can only capture linear relationship (See Pearson_wiki.png in the same folder)

--> For, non linear we use Spearman (See Spearman_wiki.png in the same folder)

X	Y	0.6
X	Z	0.7

① Spearman Rank Correlation

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

Denoted by r suffix

R(x)=rank of x
 R(y)=rank of y
 Cov = covariance
 sigma = standard deviation

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

Note that using covariance we were not able to determine strength of correlation (only direction of correlation were able to determine). So, the concept of correlation was introduced to determine strength of correlation (Pearson and Spearman). Question is that why we are doing this?

During EDA/feature engineering we need to select features that are highly correlated with output features and drop less correlated ~0. Here these concepts will help in determining the same.

Feature Selection

+ve
Size of house

+ve
No. of rooms

+ve
location
+ve correlated with o/p

≈ 0
~~No. of people Staying~~

-ve
banned

O/P
Price ↑
-ve correlated with o/p

this will have negligible impact over price of the house hence correlation is near to 0

Probability Distribution Function And Probability Density Function

Probability Mass function

Probability Distribution Function

(Below are the types of Probability Distribution Function)

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

Never consider probability distribution fun^ similar to probability density fun^.

Probability density fun^ is of the type of Probability distribution function.

When we want to draw histogram to see the distribution of data consisting of discrete random variables, we will use probability mass fun^.

① PMF ↴

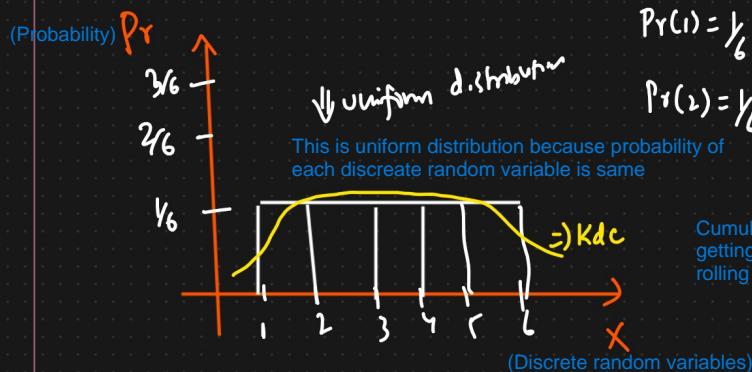
PMF characterizes the distribution of a discrete random variable

① Discrete Random Variable ↴

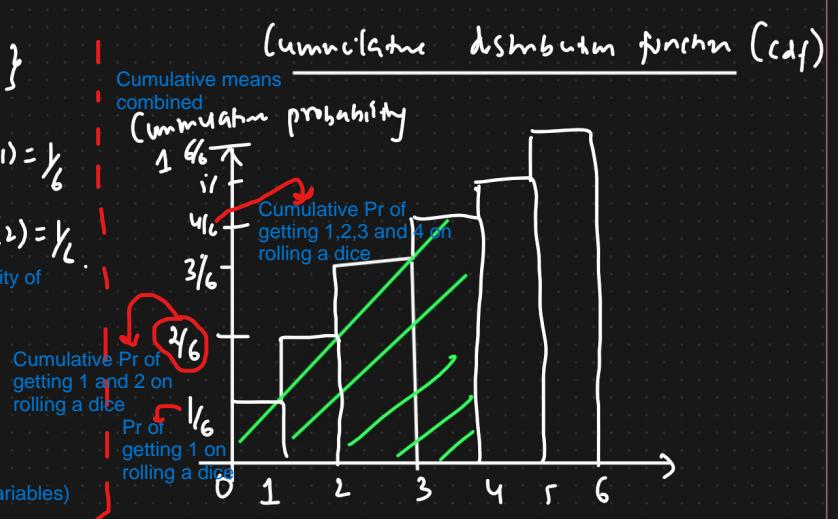
Discrete random variable:

Whole numbers with some range of values. Note that such variable holds outcomes of random experiment for eg: rolling a dice.

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



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



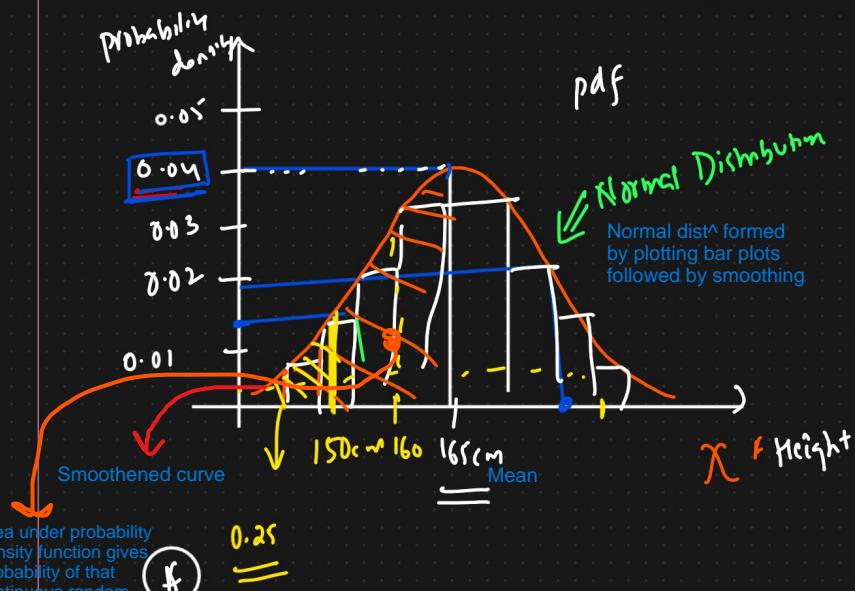
$$\Pr(X \leq 4) = \Pr(X=1) + \Pr(X=2) + \Pr(X=3) + \Pr(X=4)$$

$$= 0.67$$

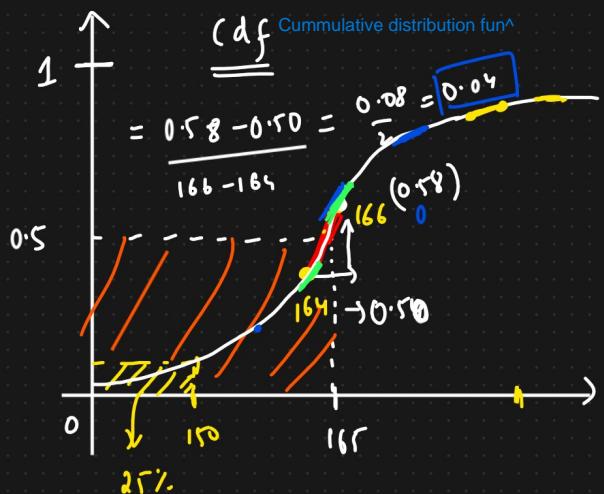
② Probability Density Function

Probability density function deals with the distribution of Continuous random variable. For eg: I am 175 cm long.

① Distribution of continuous Random Variable



Area under probability density function gives Probability of that continuous random variable. For Eg: normal distribution with mean height 165 cm, since this is mean represent 50% of distribution. So, area under this point will be 0.5 which is also the probability of occurrence of mean data point 165cm (This example is represented in above PDF and CDF using orange shaded area/section/region).

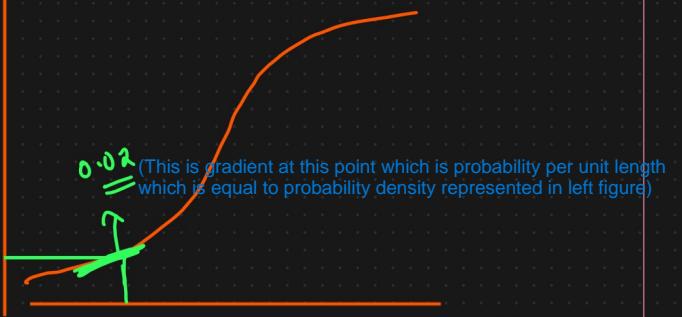
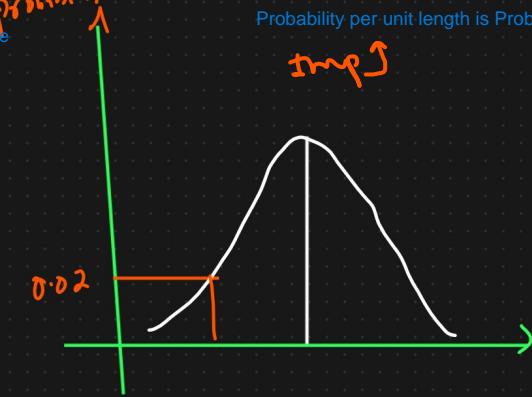


In CDF along y axis probability is taken. For normally distributed PDF the CDF thus formed is S shaped ranging from 0 to 1 on y axis (range of probability is [0,1]). Mean height is 165 that's why its projection on x axis is 0.5

Probability Density \Rightarrow Gradient of Cumulative Curve

Every probability density that we are getting in PDF is coming from gradient (which is basically slope) of respective data points in CDF.

Probability per unit length is Probability density (This statement is similar to what has been written above related to gradient)



Different types of Distribution

Why are we learning about these distributions?

Ans: because when we will be plotting our data it will follow one of these distributions. So knowing about these distributions in advance will help us understand our data better.

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

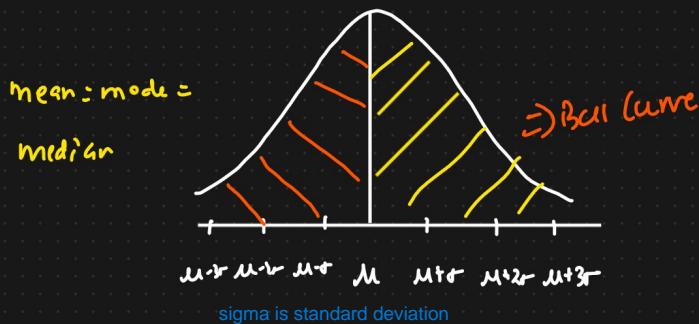
⑥ Binomial Distribution \rightarrow pmf

⑦ Poisson Distribution \rightarrow pmf

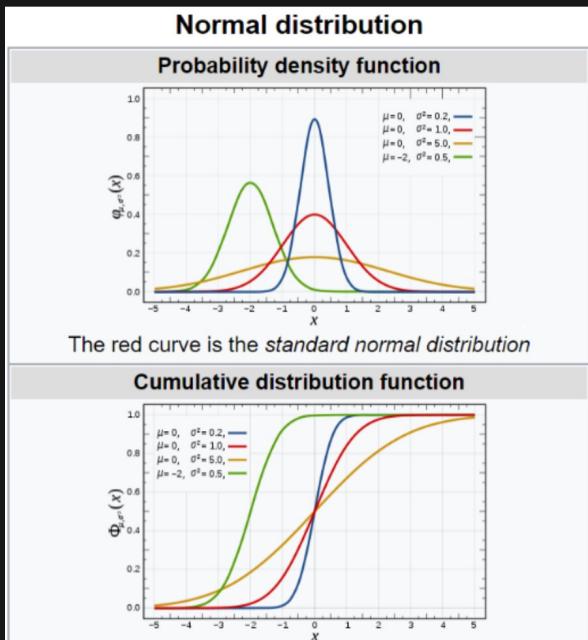
- ⑧ Uniform Distribution \rightarrow Discrete \rightarrow pmf
 \rightarrow Continuous \rightarrow pdf

- ⑨ Exponential Distribution. \rightarrow pdf
 - ⑩ CHI SQUARE Distribution \rightarrow pdf
 - ⑪ F Distribution \rightarrow pdf.

① Normal/Gaussian Distribution



Eg:- Height , weight , age , IRIS dataset



χ (⇒) Continuous Random Variable.

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

Support parameters $\mu = \text{mean}$

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

Empirical Rule

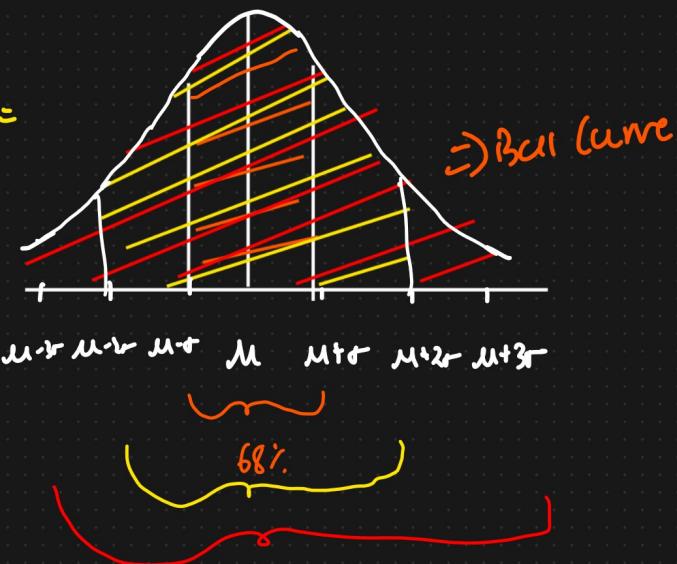
$68 - 95 - 99.7\% \text{ Rule}$

100 datapoint

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

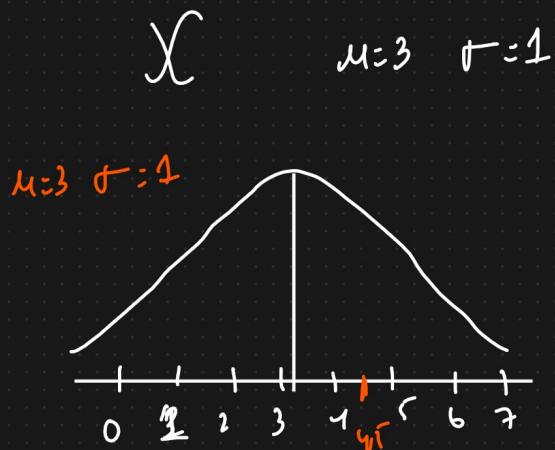
mean = mode =

median



Standard Normal Distribution

Suppose we have a random variable X distributed normally with mean=3 and standard deviation=1. If we transform it into normal distribution with mean=0 and standard deviation=1, the transformed distribution thus formed is called standard normal distribution.



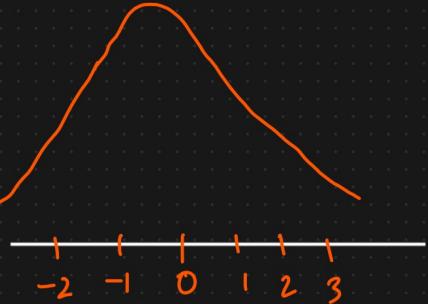
This transformation is done with Zscore

→ Standard Normal Distribution

Transformation

$$\Rightarrow$$

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



$$\downarrow$$

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

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

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

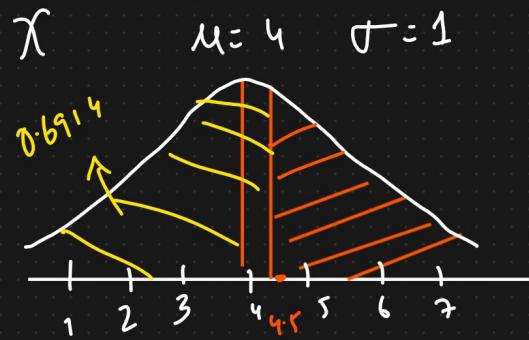
$$1 \text{ in normal dist} \Rightarrow \text{transformed to } -2 \text{ for standard normal dist}$$

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

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

Similarly, mean 3 in normal distribution get transformed into 0 in standard normal distribution.

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



11(a)

What is the percentage of scores above 6.52
 $\Rightarrow \underline{\text{Z-table}}$

$$Z\text{-score} = \frac{6.52 - 4}{1} = 2.52$$

$$\begin{aligned} \text{Area under curve} &= 1 - 0.6914 \\ &= 0.3086 // = 30.86\% \end{aligned}$$