



STATISTICS FOR THE DATA SCIENCE

Part - 2

- HISTOGRAM AND SKEWNESS
- COVARIANCE AND CORRELATION
- ADVANTAGE VS DISADVANTAGE OF VARIANCE
- PEARSON CORRELATION COEFFICIENT
- PDF / PMF / CDF
- BERNOULLY DISTRIBUTION
- BINOMIAL DISTRIBUTION

#Value_freeContent



@Krishan kumar

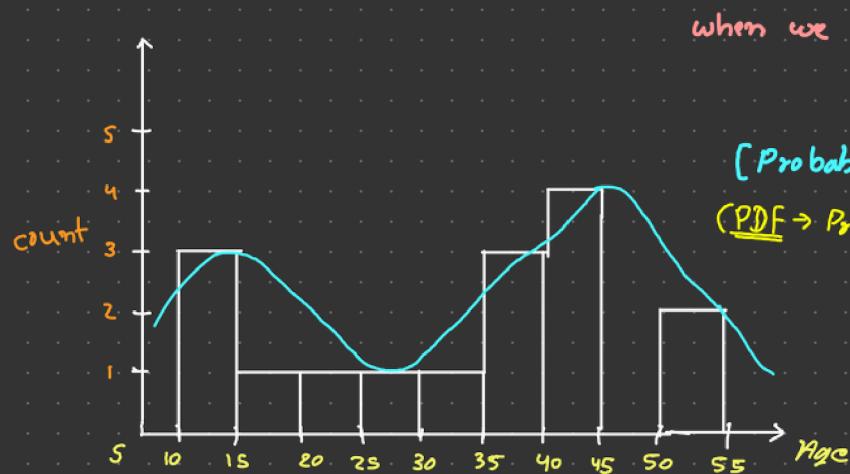
Histograms And Skewness

Age = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 45, 51}

0-50
value

$$\hookrightarrow \frac{50}{10} = 5 \text{ bin size} \quad (\text{No. of bins} \rightarrow 10)$$

$$\hookrightarrow \frac{50}{20} = 2.5 \text{ bin size} \quad (\text{No. of bins} \rightarrow 20)$$



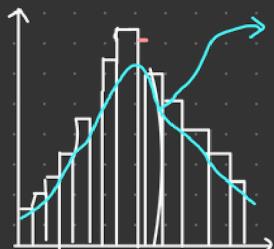
when we smoothen our histogram

↓
Technique

[Probability distribution fun.]

(PDF → Probability density function)

Skewness



Normal / Gaussian distribution

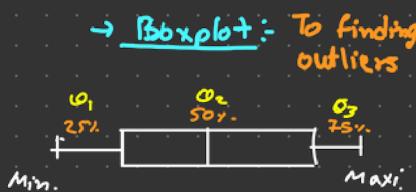
↓

These types of distribution is
called → Symmetrical distribution.

* Symmetrical Data



NO Skewness
(No curve)
equal from both
side



$$\Rightarrow Q_3 - Q_2 \approx Q_2 - Q_1$$

The mean, median, mode are all perfectly at the centre.

Mean = Median = Mode

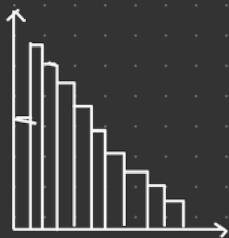


There are no skewness in the symmetrical distribution data and just because there are no any maxi value and minimum value in Right OR Left corner of data.

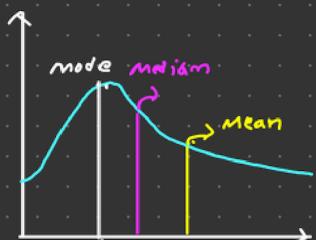
② Right Skewed data

Means

longitude in Right side (Majority data on right side present)



\Rightarrow Positive / Right Skewed \rightarrow Distribution in PDF



Relation btw
Mean, median \Rightarrow
Mode

Mean > Median > Mode

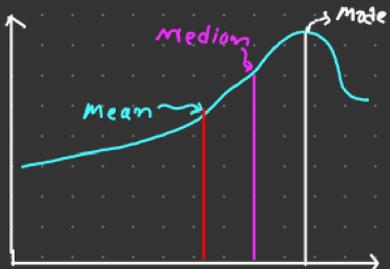


* Box plot for Right Skewed data



$$\Rightarrow Q_3 - Q_2 \geq Q_2 - Q_1$$

* Left Skewed data



* Box plot



$$\Rightarrow Q_2 - Q_1 \geq Q_3 - Q_2$$

Relation

Mean < Median < Mode



Covariance and Correlation

[To reactify Relation btw X and Y]

X	Y
2	3
4	5
6	7
8	9

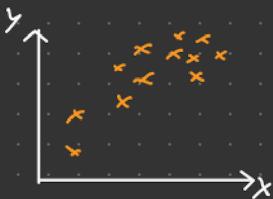
In this scenario →

X↑	Y↑
X↓	Y↑
X↑	Y↓
X↓	Y↓

For example

$x \uparrow$	$y \uparrow$
$x \downarrow$	$y \downarrow$

Plotted like →



Covariance →

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

$x_i \rightarrow$ Data point of x

$\bar{x} \rightarrow$ Sample mean of x

$y_i \rightarrow$ Data point of y

$\bar{y} \rightarrow$ Sample mean of y

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

Q What are the difference between covariance and variance?

A Varience of x $\text{var}(x)$ is nothing but $\text{cov}(x, x)$. Whenever we talk about $\text{variance}(x)$ it is specifically told about spread of the data.

⇒

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

$x \uparrow$	$y \uparrow$
$x \downarrow$	$y \downarrow$

→ Then the output will be
+ve covariance

$x \uparrow$	$y \downarrow$
$x \downarrow$	$y \uparrow$

→ Then output will be
-ve covariance

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

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

$$\Rightarrow \frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{3-1}$$

$$\Rightarrow \frac{4+0+4}{2} = 4 \rightarrow +ve \text{ value}$$

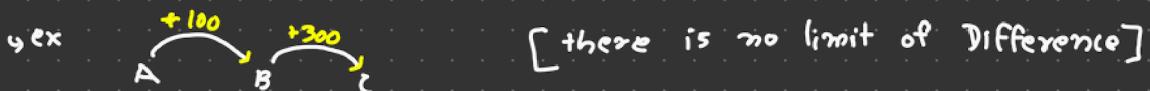
positive covariance = +ve value

So, X and Y are having a positive covariance.

Advantage of covariance,

- * Relation between X and Y +ve OR -ve value of covariance
- * Disadvantage of covariance,
- * We don't conclude weather which value is covariance with any other value Because \downarrow

Covariance does not have a specific limit value.



Note:- To fix the disadvantage of covariance we use:-
↳ Pearson Correlation coefficient

Pearson correlation coefficient

limitation $\rightarrow +1$ to -1

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

- * The more value towards $+1$ the more +ve correlated it is.
- * The more value towards -1 the more -ve correlated it is.

Spearman Rank correlation

Rank $\rightarrow -1$ to $+1$

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \times \sigma(R(y))}$$

Rank \rightarrow Highest element in the values
makes \rightarrow 1st Rank

X	Y	R(x)	R(y)
1	2	5	5
3	4	4	4
5	6	3	3
7	8	2	1
0	7	6	2
8	1	1	6

* Feature selection

+ve	+ve	+ve	= 0	-ve
Size of House	No. of Rooms	Location	No. of people staying	Houmte \Rightarrow Price $\uparrow\uparrow$

PDF / PMF

* PDF = Probability Density Function

* PMF = Probability mass function.

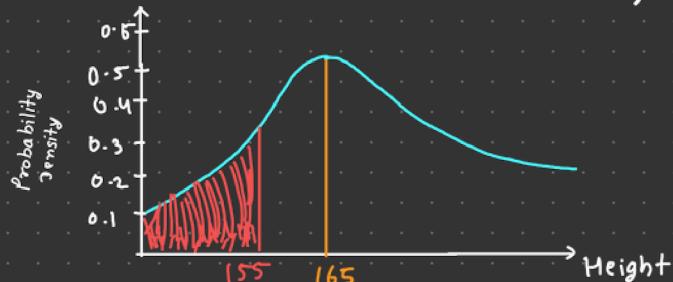
↳ PDF is denote a distribution of data. That helps to understand that how our data will distributed.

► Types of Probability Distribution function

- ① Pdf
- ② Pmf
- ③ Cdf

► Probability density function (Pdf)

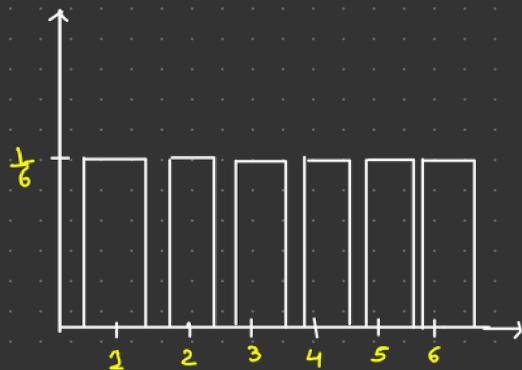
↳ Continuous Random variables, ex, Height of students in classroom



► Probability mass function (Pmf)

↪ Variable → Discrete Random variable

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



$$P(X \leq 4)$$

↓

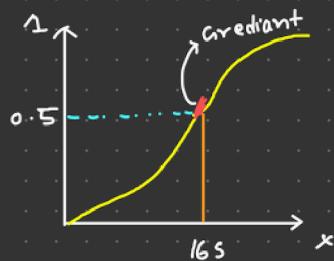
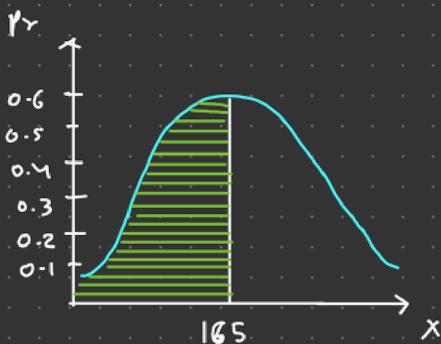
$$P(X=1) + P(X=2) + P(X=3) \\ + P(X=4)$$

$$\Rightarrow \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \rightarrow \frac{4}{6} = \frac{2}{3}$$

$$P(1) \rightarrow \frac{1}{6}$$

$$P(2) \rightarrow \frac{1}{6}$$

► Cumulative Distribution fun. (cdf),



PDF Vs PMF Vs CDF

Relation and difference

① PMF

CDF

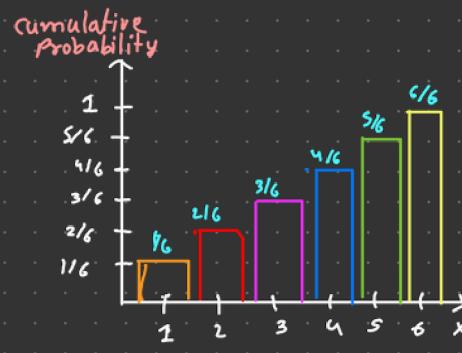
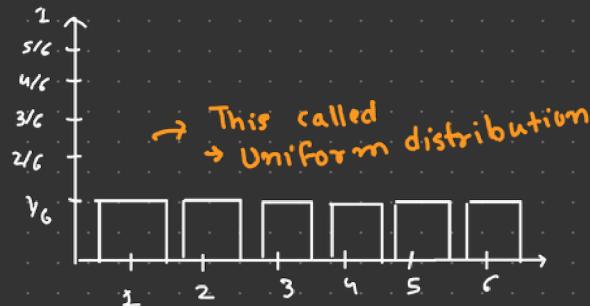
* Discrete Random Variable.

Ex → Rolling a dice $\Rightarrow [1, 2, 3, 4, 5, 6]$

$$P(1) \rightarrow \frac{1}{6}$$

⋮

$$P(6) \rightarrow \frac{1}{6}$$



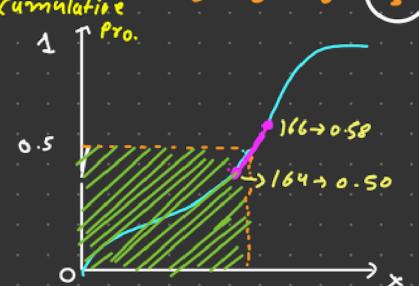
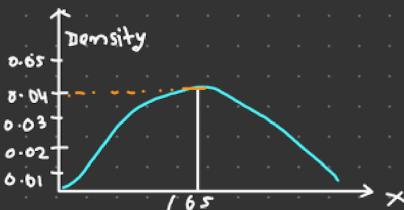
Ex → cdf for $P(X \leq 2)$

$$\text{CDF} = P(X \leq 2) = P(X=1) + P(X=2)$$

$$= \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

② PDF

↪ Distribution of continuous Random variable



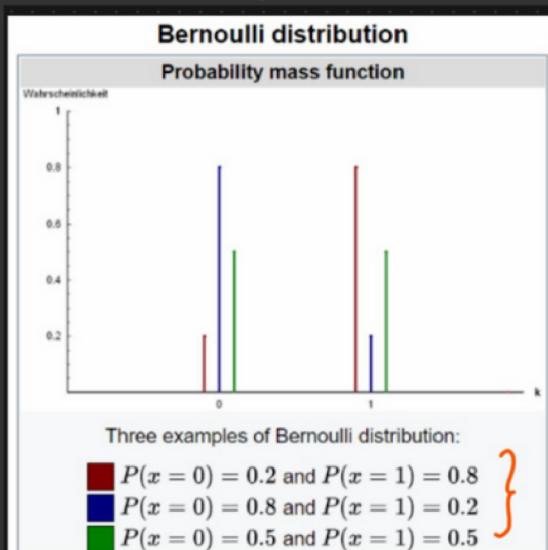
$$\Rightarrow \text{slope} \Rightarrow \frac{0.58 - 0.5}{1.66 - 1.64} \Rightarrow \frac{0.08}{0.02} \Rightarrow 0.04$$

Types of Probability Distribution

- Normal / Gaussian distribution (Pmf)
- Bernoulli distribution (Pmf)
- Uniform distribution (Pmf)
- Poisson distribution (Pmf)
- Log normal distribution (Pmf)
- Binomial distribution (Pmf)

► Bernoulli distribution → In the discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q=1-p$. Less formally it can be thought us a model for the set of possible outcomes of any single experiment that asks a yes-no question.

Pmf → To see the distribution



★ Points

- ① Discrete random variable (Pmf)
- ② Outcomes are binary
 - ex. \cong Tossing a coin (H,T)
 $P(H) \rightarrow 0.5 \rightarrow p$ (Pr. of Success)
 $P(T) \rightarrow 0.5 \rightarrow 1-p \Rightarrow q$
 $(q = \text{Pr. of Failure})$
- ③ When there will fail / pass
 - $P(\text{Pass}) \rightarrow 0.7 \rightarrow p$ Success value
 - $P(\text{Fail}) \rightarrow 1-0.7 \Rightarrow 0.3 = q$ Failure value

$$* \underline{\text{PMF}} = P^k * (1-P)^{1-k}, \quad k \in \{0,1\}$$

if $k=1$

$$\text{PMF} \rightarrow P(k=1) = P^1 * (1-P)^{1-1}$$

$$= P$$

if $k=0$

$$\begin{aligned} \text{PMF} &= P(k=0) = P^0 * (1-P)^{1-0} \\ &= (1-P) = q \\ &\Rightarrow q \end{aligned}$$

Simplified

$$\text{PMF} = \begin{cases} 1-P & , \text{ if } k=0 \\ P & , \text{ if } k=1 \end{cases}$$

\Rightarrow Mean of Bernoulli distribution,

$$\langle \varepsilon(k) \rangle = \sum_{j=1}^k k \cdot P(k)$$

$$, \quad P(k=1) = 0.6 \Rightarrow P$$

$$, \quad P(k=0) = 0.4 \Rightarrow 1-P \Rightarrow q$$

$$\Rightarrow [0 \times 0.4 + 1 \times 0.6] \Rightarrow 0.6 = P$$

Whenever we abstract mean in Bernoulli distribution, then we get P-Value.

\Rightarrow Median of Bernoulli dis,

$$\text{Median} = \begin{cases} 0 & , \text{ if } P < \frac{1}{2} \\ [0,1] & , \text{ if } P = \frac{1}{2} \\ 1 & , \text{ if } P > \frac{1}{2} \end{cases}$$

\Rightarrow Variance of Bernoulli dis.,

$$\text{Var} = P(1-P)$$

$$\boxed{\text{Var} = pq}$$

\Rightarrow Std of Bernoulli dis.,

$$\boxed{\text{Std} = \sqrt{pq}}$$

► **Binomial Distribution** In Statistic the binomial dis. with parameter n and P is the discrete P.d. Dis of the number of success in a sequence of n independent experiments. Each asking Yes or No. question. and each its own boolean-value outcome. success (with Pro. P) or failure ($q=1-p$).

A single success / failure experiment is also called a bernoulli trial OR Bernoulli experiment and a sequence of outcome is called a bernoulli process.

For a single trial ex, $n=1$ the binomial dis. is a bernoulli dis. the binomial dis. is the basis for the popular binomial test of statistical significance.

$$B(n,p)$$

Notation: $B(n,p)$

- * Every experiment outcome is binary
- * This experiment is performed for n trials
- * Group of bernoulli dis \rightarrow Binomial Distribution

Parameters

- * $n \in \{0, 1, 2, \dots\} \rightarrow$ No. of trials
- * $p \in [0, 1] \rightarrow$ Success Pro. of each trial
- * $q = 1-p$

* Ex Tossing a coin 10 times

Support :- $k \in \{0, 1, 2, 3, \dots, n\}$
 \rightarrow Number of success

* Pmf = $P(k, n, p) = {}^n C_k p^k (1-p)^{n-k}$

for $k = 0, 1, 2, 3, \dots, n$ where

$${}^n C_k = \frac{n!}{k! (n-k)!}$$

\Rightarrow Mean :-

$$\text{Mean} = np$$

\Rightarrow Variance :-

$$\text{Var} = npq$$

\Rightarrow Std

$$\text{Std} = \sqrt{npq}$$