

## Agenda

→ Covariance & Correlation

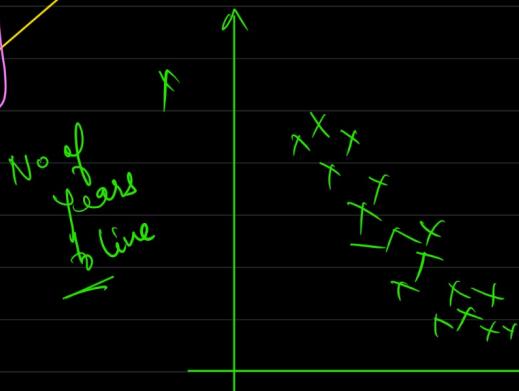
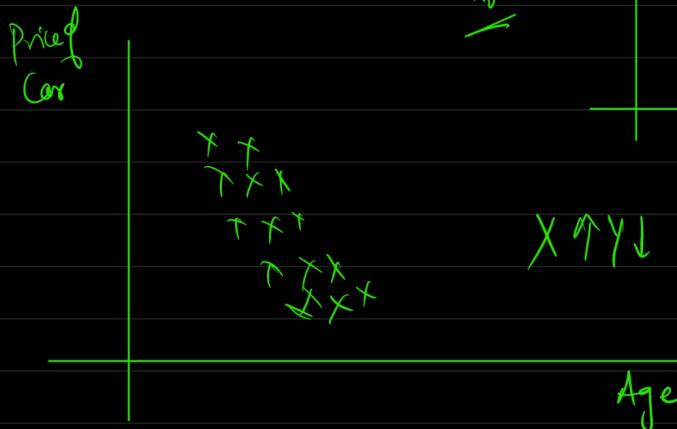
→ Probability dist<sup>n</sup>

Till now  
→ Descriptive Stats

→ Date

	Area of house	No of room	Price of house
→	1000	2	19.
→	-	-	-
→	-	-	-
→	-	-	-
→	-	-	-
→	-	-	-

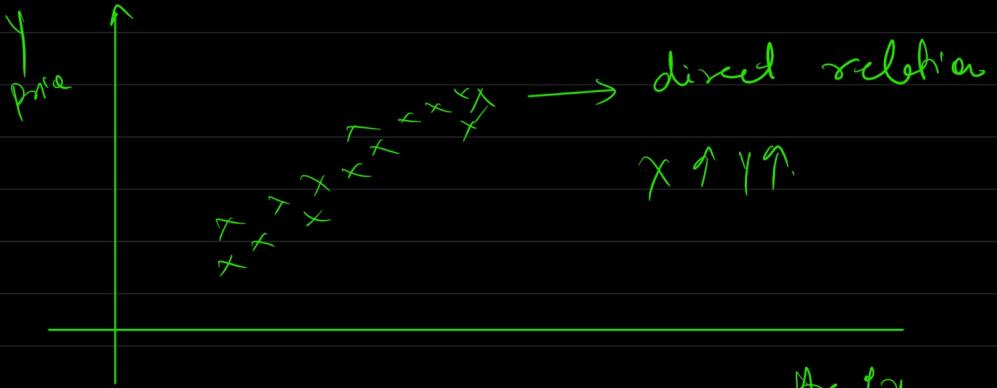
Area of House ↑ Price of house ↑



$X \uparrow Y \downarrow$   
(indirect)

Absolute consumption

X



$X \uparrow Y \uparrow$

Average

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

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

Variance was spread of data  $\Rightarrow$  relationship of feature with itself

$\text{Cov} \rightarrow$  w.r.t other featu

Sample

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

$\star \text{Cov} \rightarrow$  relationship b/w two variable

$\left\{ \begin{array}{l} X \uparrow Y \uparrow \\ X \downarrow Y \downarrow \end{array} \right.$   
 $(+ve \text{ cov})$

$X \uparrow Y \downarrow$   
 $X \downarrow Y \uparrow$

$(-ve \text{ cov})$

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

S.no	X	Y
1	2	5
2	3	2
3	4	1
4	5	6

$$\Rightarrow (2-2.75)(5-4) + (3-2.75)(2-4) + (4-2.75)(1-4)$$

$$+ (5-2.75)(6-4)$$

$$\bar{x} = 2.75 \quad \bar{y} = \frac{16}{4} = 4$$

$$4-1$$

$$\Rightarrow -46.33$$

$\curvearrowleft$   $X$  and  $Y$   $\rightarrow$   $\text{vc}$

$\Rightarrow$  +ve value  $\rightarrow$  +ve relationship

Advantage

$\rightarrow$  Relation b/w  $X$  &  $Y$   $\rightarrow$  +ve,  $\underline{-\text{ve}}$

\* Disadvantage

①  $X$  .  $Y$

$$\text{Cov}(X, Y) = 50$$

$\downarrow$   
 $\rightarrow \infty \rightarrow \infty$

$A - B$

$$\text{Cov}(A, B) = 100.$$

$A - B$  is twice the relationship as compared to  $X - Y$ .

$\rightarrow$  Covariance only measures the direction.

$\rightarrow$  Never tells about Strength of relation

$\rightarrow$  No any standardise scale to interpret the strength.





$\text{Cov} \rightarrow$  direction of relation



② Covariance has dimension

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

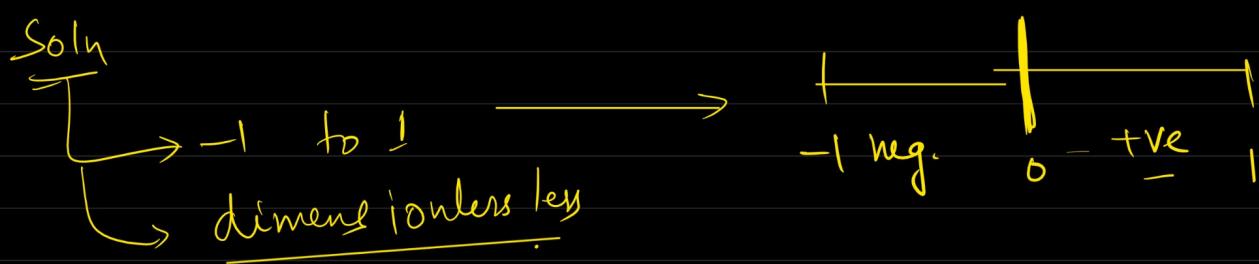
$x \rightarrow \text{height} \rightarrow \text{ft}$   
 $y - \text{wt} \rightarrow \text{kg}$   $\text{Cov}(x, y) \Rightarrow \underline{\text{ft} \cdot \text{kg}}$

$$\text{Cov}(\text{transient}, \text{ht}) = \underline{\underline{R_s}} \cdot \underline{\underline{\text{ft}}}.$$

$x$      $y$      $z$   
 $\downarrow$      $\swarrow$      $\searrow$   
 $\text{ft} \cdot \text{kg}$      $\text{ft} \cdot R_s$

$f(x)$ (Area)	$\text{Y}_{\text{avg}}$	$P_{\text{rec}}$

Not comparable  
 $\hookrightarrow$  different dimensions



$$\frac{f_1 f_2}{f_3 f_4}$$

② Pearson Correlation Coeff [-1 to 1]

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



$$\xleftarrow{-\nu_c} \quad \xrightarrow{+\nu}$$

the more +ve  
relationship is  
higher on  
+ve direction

the more -ve  
correlation is  
more on negside

$$\rho_{x,y} = 0.4, \quad \rho_{A,B} = 0.8.$$

$\Rightarrow$  feature A, B is highly correlated as compared to x, y.

\* → Pearson correlation coefficient always measures the linear relationship.

What to do in Non-linear relationship ??

\* Spearman Rank Correlation.

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

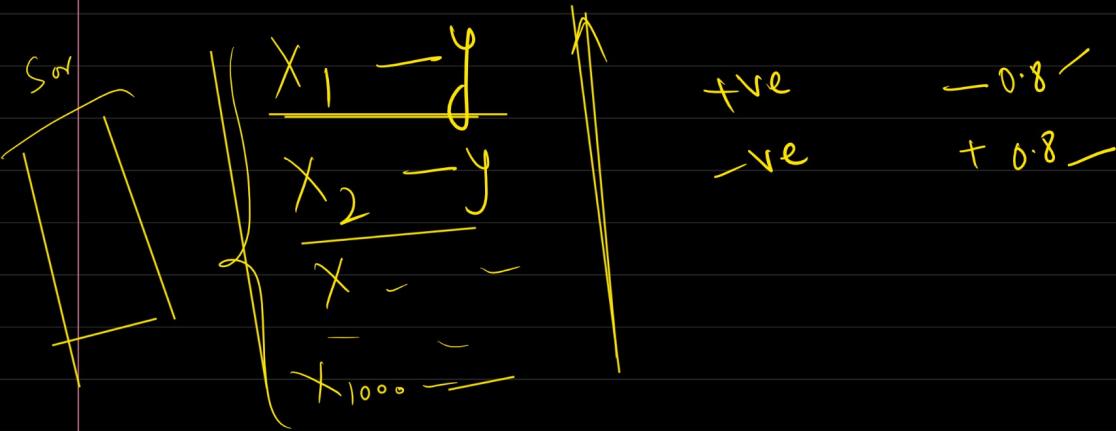
$\downarrow$   
 $x$   $R(x)$   $y$   $R(y)$

5	3	6	1
7	2	4	2
8	1	3	3
1	5	1	5
2	4	2	4

$X \rightarrow 5, 7, 8, 1, 2 \rightarrow 1, 2, 3, 4, 5$

$R(x) \rightarrow$  Rank of  $x$   
 $R(y) \rightarrow$  Rank of  $y$

$x_1, x_2, x_3, \dots, x_{1000} \xrightarrow{\text{rank}} p_{\text{rinc}}$



## \* Random Variables

↳ A set of possible values from a random experiment.

→ Tossing a coin



$$X = \begin{cases} 0 : - \text{ tail} \\ 1 : - \text{ Head} \end{cases}$$

random Variable

quantifies the experiment

$$\lambda = \{1, 0\}$$

$$P(X) = \frac{1}{n} \quad \left( \begin{array}{l} \text{where } n \text{ is total possible} \\ \text{outcomes} \end{array} \right)$$

\* dice

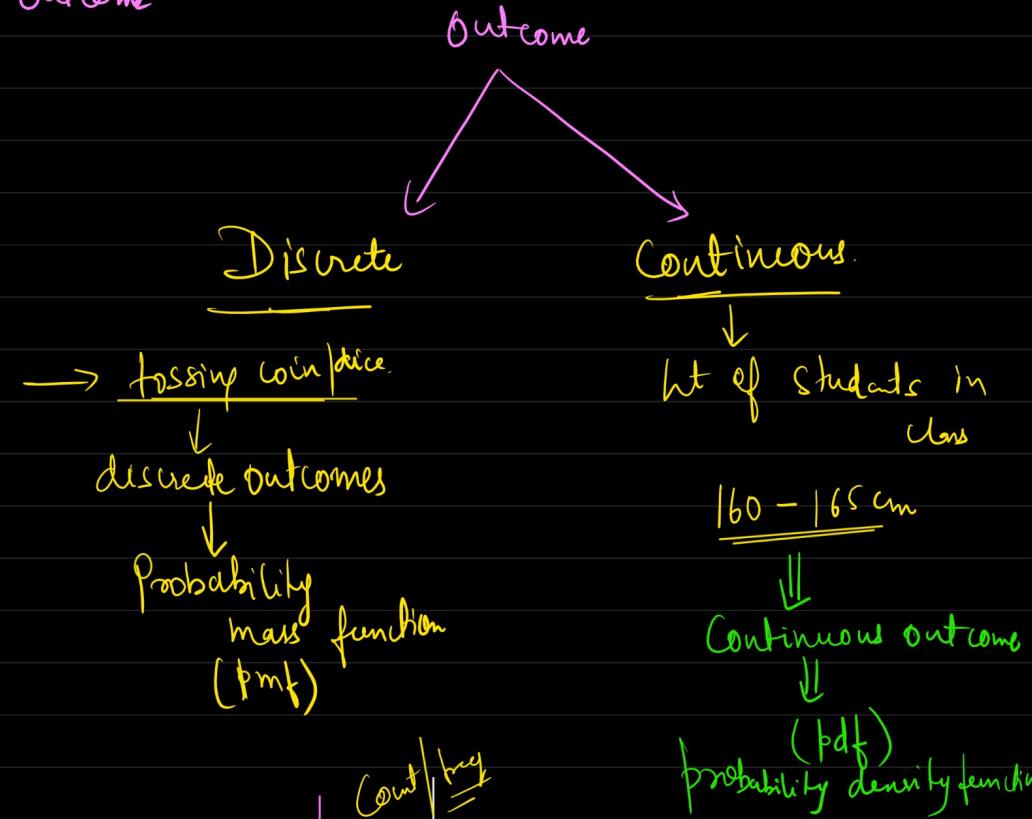
→ 1, 2, 3, 4, 5, 6

$$\begin{matrix} \downarrow & \downarrow & \downarrow & & & \downarrow \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \cdot & \cdot & \frac{1}{6} \end{matrix}$$

$$\boxed{P(X) = \frac{1}{n}} \quad (\text{where } n=6)$$

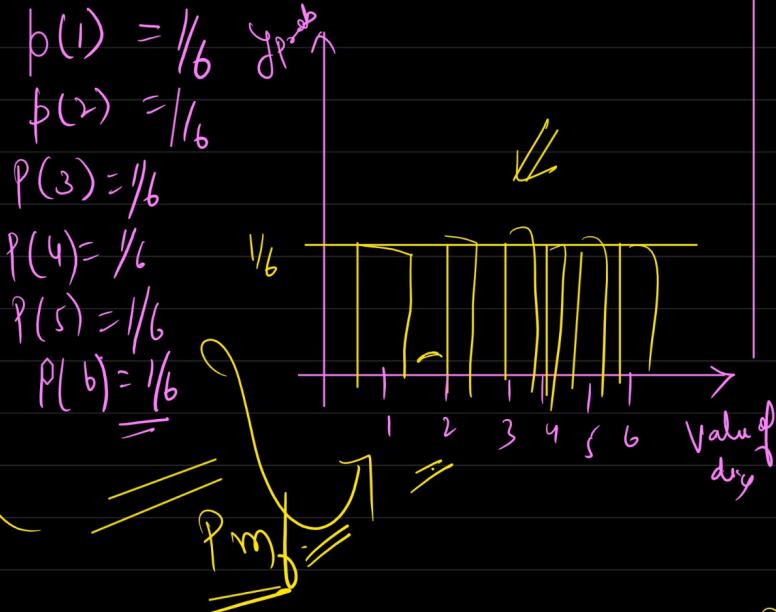
$$\frac{P(X=1)}{P(X=2)} = \frac{1}{6} = \underline{\underline{\frac{1}{6}}}$$

Dice - 1, 2, 3, 4, 5, 6  
Outcome

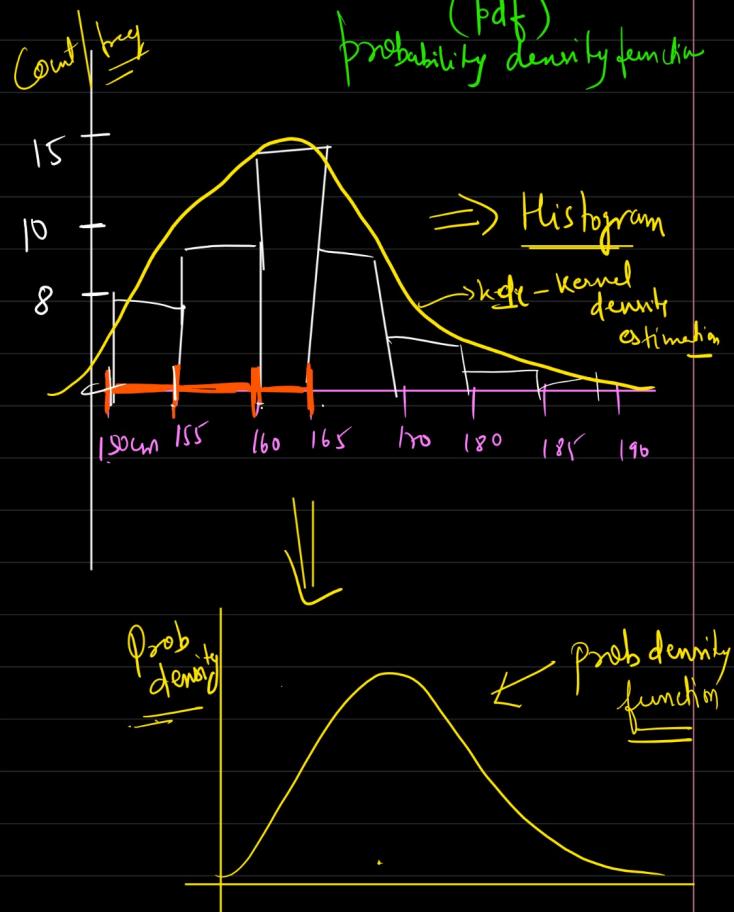


### Discrete

Dice → Discrete outcomes



Categorical data - Bar  
 Continuous - Historical



Discrete — 1, 2, 3, 4 (which takes finite value)

Continuous → ht of stn — 1, 2  
 150 cm,  
 150.001 cm  
 100 kg  
 100.2 pairs

Two types of Experiment

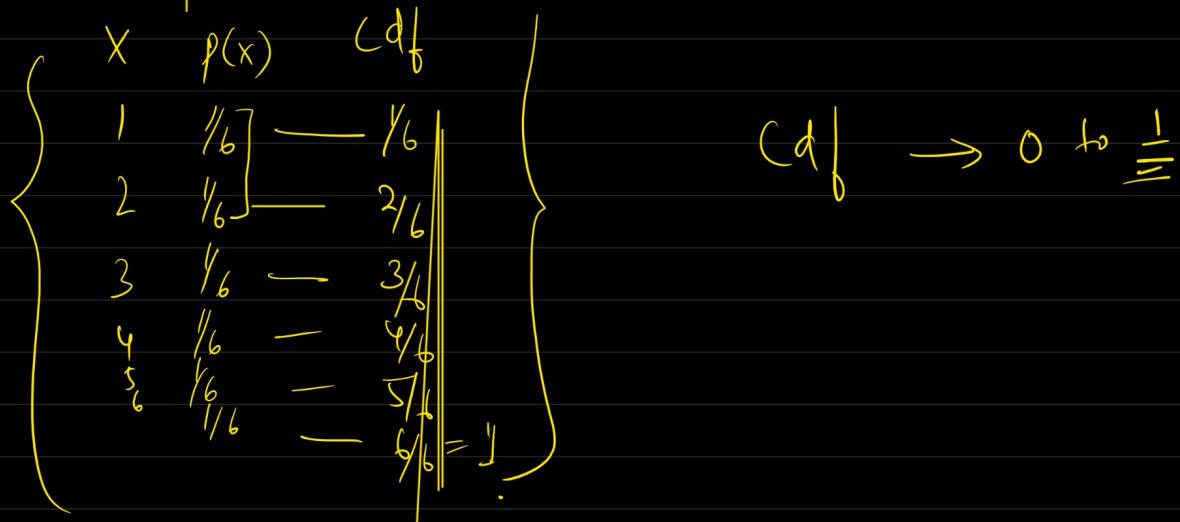
discrete

continuous

↓  
Prob mass function

↓  
Prob density function

Cumulative distn function



## Prob dist fn

### Prob density fn

### Prob mass fn.

→ Normal / Gaussian → Bernoulli

→ Standard Normal → Binomial

→ log Normal dist → Poisson

→ Chi-square dist

→ F-dist.

discrete  
Uniform dis cont. prob

### \* Discrete Uniform dist<sup>h</sup>

→ Outcomes will be discrete

→ Outcomes are equally likely.

### Uniform dist

Discrete  
(pmf)

Continuous  
(pdf)

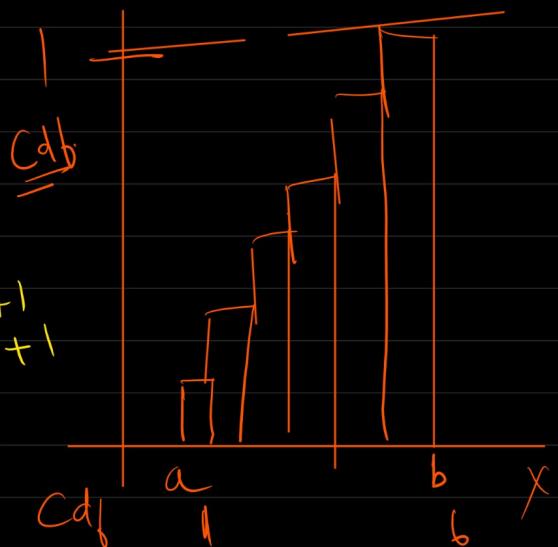
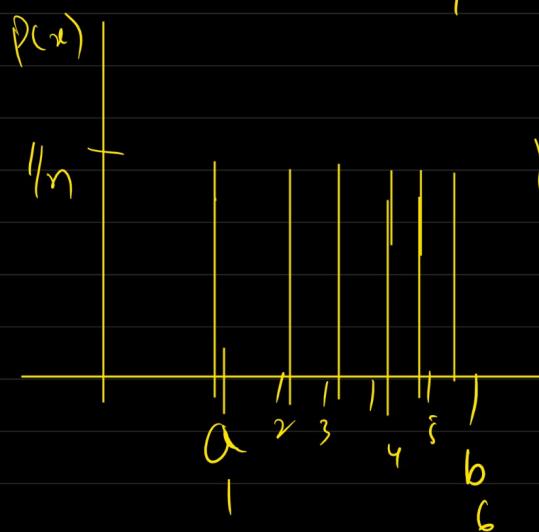
→ Outcomes are discrete &  
equally likely

Same/Equally probability

of rolling a dice

$$\begin{aligned} P(1) &= \frac{1}{6} \\ P(2) &= \frac{1}{6} \\ &\vdots \\ P(6) &= \frac{1}{6} \end{aligned}$$

$$U(a, b) \mid \text{Unif } \{a, b\}$$



Q What is prob of getting 3 when you throw a dice?

$$\rightarrow \frac{1}{6}$$

\* Mean of discrete Uniform dist  $\Rightarrow \frac{a+b}{2}$  ✓

Variance "  $\Rightarrow \frac{n^2-1}{12}$  ✓

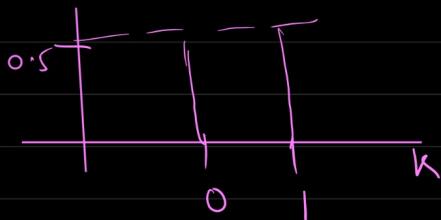
## ② Bernoulli dist<sup>n</sup>

→ discrete  
→ Outcomes — 2 → Pass | fail.

→ Tossing a coin.

$$P(X=k) = \begin{cases} p & \text{if } k=1 \\ 1-p & \text{if } k=0 \end{cases}$$

$$P(X=k) = p^k (1-p)^{1-k}$$



$$\rightarrow P_{\text{Success}} + P_{\text{failure}} = 1$$

$$\rightarrow P_{\text{Head}} + P_{\text{tail}} = 1$$

$$\rightarrow 0.5 + 0.5 = 1$$

$$P(x=k) = p^k (1-p)^{1-k}$$

$$k=1$$

$$k=0$$

$$p^1 (1-p)^{1-1}$$

$\cancel{p^1}$

$$(1-p)^0$$

$\cancel{(1-p)}$

$p$

$$p^0 (1-p)^{1-0}$$

$\downarrow$

$\downarrow$

$(1-p)$

$$p + (1-p) = 1$$

\* Conditions of Bernoulli dist

- No of trials should be finite
- Each trial should be independent.
- Only two possible outcomes
- Prob of each outcome should be same in every trial.

Example:

tossing a coin

head/tail

yes/no

Pass/fail

Bumrah bowls 6 balls at wicket with prob of 0.6 at hitting the stump with each ball. What is prob not hitting wicket

$$p(H) = 0.6$$

$$p(\text{Not } H) = 1 - 0.6 = 0.4$$

$$\begin{cases} \text{mean} = p \\ \text{variance} = p(1-p) \end{cases}$$

\* Bi-nomial dist



bi-two  $\rightarrow$  bernoulli

n-bernoulli  $\rightarrow$  Binomial

\* Binomial dist<sup>n</sup> is n Bernoulli trial

$$\text{pmf of Bernoulli} = p^k (1-p)^{1-k}$$

$$\text{pmf of Binomial dist} = {}^n C_k p^k (1-p)^{n-k}$$



1st toss 2nd toss 3rd toss 4th, ..., 10<sup>th</sup>

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

\* What is prob of getting k heads out of n trials

n  $\rightarrow$  total no. of trials  
k - no of events.

$$\begin{aligned} 10 &= 4 \times 3 \times 2 \times 1 \\ 10 &- 10 \times 9 \times 8 \times \dots \times 1 \end{aligned}$$

Q With 3 tosses what is prob of getting exactly 2 heads?

$$n=3, k=2$$

$$p(k=2) = {}^n C_k p^k (1-p)^{n-k}$$

$$= {}^3 C_2 (0.5)^2 (0.5)^{3-2}$$

$$\text{mean} = np$$

$$\text{Variance} = np(1-p)$$

$$= \frac{3}{3-2 \cdot 1} (0.5)^2 \cdot 0.5$$

$$= \frac{3 \times 1}{1 \times 2} (0.5)^3 = 0.375$$

### \* Poisson distribution

→ discrete Uniforms.

→ describes the no. of events that occur within a fixed interval of time, given average rate of occurrence.

→ No of events occurring in a fixed time interval.

$$\text{pmf} \rightarrow p(x=k) = \frac{-\lambda^k}{k!} e^{-\lambda}$$

$$e = 2.718$$

$\lambda$  - avg rate of events every interval.

- ex No of calls received by a customer care every hour  
 ex No of People visiting hospital / bank  
 ex No of accidents  
 ex No of emails received

Q The avg no of customers entering a store in an hour is 5. What is the prob of exactly 3 customers will enter the store next hour?

$$\lambda = 5$$

$$P(X=3)$$

$$\frac{e^{-\lambda} \lambda^x}{x!} = e^{-5} \frac{5^3}{3!}$$

$$= \underbrace{(2.718)^{-5}}_6 \times 125$$

$\text{Mean} = \lambda t$
$\text{Varian} = \lambda t$

$$= \approx 0.14$$