

Discrete Probability Distributions

Infinite trials:

ex:- let's see who wins?

lets consider A & B and both throw the dice.

whoever gets the first six wins. Consider

A starts.

$$P(A \text{ winning}) = \frac{1}{6} + \frac{5}{6} \times \frac{5}{6} \times \frac{1}{6} + \frac{5}{6} \times \frac{5}{6} \times \frac{5}{6} \times \frac{5}{6} \times \frac{1}{6} + \dots$$

in first chance in second chance ... up to ∞

$$= \frac{1}{6} + \left(\frac{5}{6}\right)^2 \times \frac{1}{6} + \left(\frac{5}{6}\right)^4 \times \frac{1}{6} + \dots$$



Infinite geometric progression.

$$\therefore S_{\infty} (\text{Sum of infinite terms}) = \frac{a}{1 - R}$$

$a \rightarrow$ first term

$R \rightarrow$ common Ratio

here

$$a = \frac{1}{6} \quad R = \left(\frac{5}{6}\right)^2 \rightarrow \frac{\text{3rd term divided by 2nd term}}{}$$

$$P(A \text{ winning}) = \frac{\frac{1}{6}}{1 - \left(\frac{5}{6}\right)^2} = \frac{\frac{1}{6}}{1 - \frac{25}{36}}$$

$$= \frac{6}{11}$$

So whoever starts 1st will have $\frac{6}{11}$ probability to win.

So person starting 2nd will have $\frac{5}{11}$ probability to win. (i.e $1 - \frac{6}{11} = \frac{5}{11}$)

(or)

$$P(B \text{ winning}) = \frac{5}{6} \times \frac{1}{6} + \frac{5}{6} \times \frac{5}{6} \times \frac{5}{6} \times \frac{1}{6} + \dots$$

$$a = \frac{5}{36}$$

$$R = \left(\frac{5}{6}\right)^2$$

$$P(B \text{ winning}) = \frac{\frac{5}{36}}{1 - \frac{25}{36}} = \frac{5}{11}$$

Q) A and B throw a pair of dice. A wins if he throws 6 before B throws 7 & B wins if he throws 7 before A throws 6. If A begins find the probability of

- (i) A winning (ii) B winning.

$$n(S) = 36$$

$$n(\text{total}=6) = 5$$

$$n(\text{not } 6) = 31$$

$$n(\text{total}=7) = 6$$

$$n(\text{not } 7) = 30$$

$$P(A \text{ winning}) = \frac{5}{36} + \frac{31}{36} \times \frac{30}{36} \times \frac{5}{36} + \left(\frac{31}{36} \times \frac{30}{36} \right)^2 \times \frac{5}{36} \\ + \dots$$

$$= \frac{\frac{5}{36}}{1 - \frac{31}{36} \times \frac{30}{36}} = \frac{\frac{5}{36}}{1 - \frac{930}{1296}}$$

$$= \frac{30}{61}$$

$$P(B \text{ winning}) = 1 - P(A \text{ winning})$$

$$= \frac{31}{61}$$

$P(B \text{ winning}) > P(A \text{ winning})$ because probability of total = 7 > total = 6.

Bernoulli's Trials

- finite in number
- independent of each other
- Each trial has only 2 outcomes success
failure
- Probability of success or failure remains same in each trial.

* Bernoulli trial is also called binomial trial.

ex:- The product can be defective or not defective

Q) Eight balls are drawn from a bag containing 10 white and 10 black balls. Predict whether the trials are Bernoulli trials if the ball drawn is replaced and not replaced.

(1) with replacement the probability of success say white ball remains same so Bernoulli trials.

(2) without replacement the probability of success changes so not Bernoulli trials.

Binomial distribution:

Binomial theorem

General term $(x+y)^n = {}^n C_r x^r y^{n-r}$.

let us consider $P(\text{success}) = P$

$P(\text{failure}) = q$

We know that $P+q=1$

If we perform a trial 'n' times in which we have 'r' successes and $(n-r)$ failures.

- these 'r' successes in 'n' trials can be anywhere

so here ${}^n C_r = \frac{n!}{(n-r)! r!}$

ex:- if 10 trials, 4 successes then permutations

is ${}^{10} C_4 = \frac{10!}{6! 4!} = \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2}$

= 210

In Binomial - Distribution

$$1 = {}^n C_0 p^0 q^n + {}^n C_1 p^1 q^{n-1} + {}^n C_2 p^2 q^{n-2} + \dots$$

$\downarrow \qquad \downarrow \qquad \downarrow$

Probability of getting 0 successes in n trials	Probability of getting exactly 1 success in n trials	Probability of getting exactly 2 successes in n trials
---	--	--

$${}^n C_r p^r q^{n-r}$$
 exactly 'r' successes
in n trials

Q) An experiment succeeds twice as often as it fails. Find the probability that in next 6 trials, there will be atleast 4 successes.

$$\therefore P = 2q \quad P + q = 1$$

$$\therefore q = \frac{1}{3}, \quad P = \frac{2}{3}$$

$P(\text{atleast 4 successes}) \therefore \{4, 5, 6\}$

$$= {}^6 C_4 \left(\frac{2}{3}\right)^4 \left(\frac{1}{3}\right)^2 + {}^6 C_5 \left(\frac{2}{3}\right)^5 \left(\frac{1}{3}\right)^1 + {}^6 C_6 \left(\frac{2}{3}\right)^6 \left(\frac{1}{3}\right)^0$$

Random Variable: In probability, a real-valued function, defined over the sample space of a random experiment, is called a random variable.

Definition: A Random Variable is a rule that assigns a numerical value to each outcome in a Sample Space.

Random Variables

- Discrete Random Variable - A random Variable is said to be discrete if it assumes only Specified values in a interval

- Continuous Random Variable

If we want to create discrete Values from Continuous we can consider intervals and bins and these will be discrete

Random variable is denoted by Capital 'x'. It takes input as all the outcomes in the sample space as input and gives a number as output

Ex:- If we toss 2 coins and assign random Variable number of heads

$$S = \left\{ \begin{array}{l} HH \\ HT \\ TH \\ TT \end{array} \right\}$$

$$x - HH \xrightarrow{\text{output}} 2$$

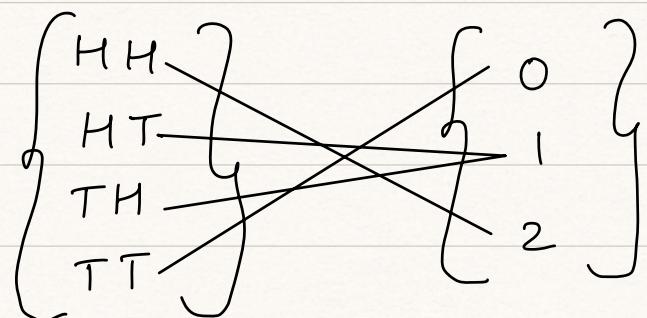
$$x - HT \rightarrow 1$$

$$x - TH \rightarrow 1$$

$$x - TT \rightarrow 0$$

so

X



* Random Variable will have many to one kind of relationship.

Discrete Random variable :

- A discrete random variable can take only a finite number of distinct values such as 0, 1, 2, ... and so on.
- The probability distribution of a random variable has a list of probabilities compared with each of its possible values known as probability Mass function.

ex:- 1) Toss a coin and count no.of heads
2) Defects in a product produced.

Continuous Random Variable:

- A numerically valued variable is said to be continuous, if the random variable X can assume an infinite and uncountable set of values, it is said to be a continuous random variable.
- when X takes any value in a given interval (a, b) , it is said to be a continuous random variable in that interval.
- A continuous random variable is such that it has constant cumulative function throughout.
- ex:- pressure, height, weight, volume etc...

Probability Distribution Function :

A function which is used to define the distribution of a probability is called a probability distribution function.

Probability Distribution Table :

A Probability distribution is a table that displays the probability that a random variable takes on certain values

X	x_1	x_2	x_3	x_n
$p(x)$	p_1	p_2	p_3	p_n

Properties of Probability Distribution table:

1) All the Probabilities must add up to 1.

$$\text{i.e } p_i \geq 0 \text{ and } p_1 + p_2 + p_3 + \dots + p_n = 1$$

2) Mean of random variable:

$$\text{In statistics mean} = \frac{\sum x_i f_i}{\sum f_i}$$

$$\text{So here mean} = \frac{\sum x_i p_i}{\sum p_i} \text{ but } \sum p_i = 1$$

$$\therefore \text{mean} = \sum x_i p_i$$

For population mean is μ

3) Variance of a Random Variable

The variance tells how much is the spread of random Variable x around the mean value.

The formula of the variance of a random variable is given by:

$$\text{Var}(x) = \sigma^2 = E(x^2) - [E(x)]^2$$

where $E(x^2) = \sum x_i^2 p_i$ and $E(x) = \sum x_i p_i$

$$\text{or } \sigma^2 = \mu(x_i^2) - (\mu(x))^2$$

here $E(x)$ is expectation of x .

$$\text{or } \sigma^2 = \sum (x_i - \mu)^2 * p(x_i)$$

x_i : The i^{th} value

μ : The mean of distribution

$p(x_i)$: The probability of i^{th} value.

4) Standard deviation is $\sqrt{\text{Variance}}$.

Q1) In the toss of 2 coins example from above.

X	0	1	2
$P(X)$	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{4}$

- The top row shows the values of Random Variable X .
- The bottom row shows the probability of the X . This is similar to frequency tables in statistics.

$$\text{mean} = 0 \times \frac{1}{4} + 1 \times \frac{2}{4} + 2 \times \frac{1}{4} = 1$$

$$\begin{aligned}\text{Variance} &= \left(0^2 \times \frac{1}{4} + 1^2 \times \frac{2}{4} + 2^2 \times \frac{1}{4}\right) - 1 \\ &= \frac{2}{4}\end{aligned}$$

Q2) A urn contains 4 white and 6 red balls.

Four balls are drawn at random from the urn. Find the probability distribution of the number of white balls.

↓ No. of white balls drawn

X	0	1	2	3	4
$P(X)$	$\frac{6C_4}{10C_4}$	$\frac{4C_1 \times 6C_3}{10C_4}$	$\frac{4C_2 \times 6C_2}{10C_4}$	$\frac{4C_3 \times 6C_1}{10C_4}$	$\frac{4C_4}{10C_4}$

Q) If a six-face die is thrown 3 times.
 what is the probability that 5 appears
 exactly 2 times.

$$n = 3$$

success is rolling a 5, failure rolling anything
 other than 5.

$$p = \frac{1}{6} \quad \text{and} \quad q = \frac{5}{6}$$

$$r = 2$$

$$P(X=x) = {}^3C_2 \left(\frac{1}{6}\right)^2 \left(1 - \frac{1}{6}\right)^{3-2}$$

$$= 0.0694$$

logically

$p^x (1-p)^{n-x} \rightarrow$ one specific ordering of one success
 and $n-x$ failures.

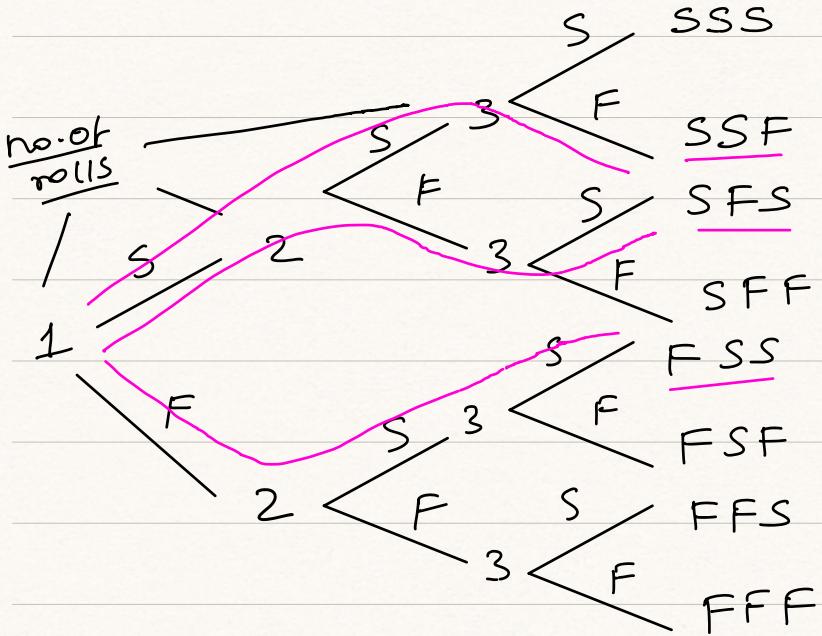
Here we are concerned about 'x' success in
 'n' trials so

$$P(X=x) = {}^nC_x p^x (1-p)^{n-x} .$$

using probability tree with the above example

$$n = 3$$

$$r = 2 \quad (2 \text{ times } 5)$$



here exactly 2 times success is (SSF, SFS, FSS)

$$\begin{aligned}
 P(X=2) &= p^2(1-p) + p(1-p)p + (1-p)p^2 \\
 &= 3p^2(1-p) \\
 &= {}^3C_2 p^2(1-p)
 \end{aligned}$$

$$\therefore P(X=x) = {}^nC_x p^x (1-p)^{n-x}$$

Mean and variance of Binomial distribution

Probability of success = P

Probability of failure = q

Total number of trials = n

Total number of successes = r

$$\text{mean} = np$$

$$\text{variance} = npq.$$

X	0	1	2	...	n
$P(X)$	$n C_0 p^0 q^n$	$n C_1 p^1 q^{n-1}$	$n C_2 p^2 q^{n-2}$	$n C_n p^n q^0$

Using binomial theorem derivation we get

$$\text{since } (p+q)^n = 1^n = 1$$

$$1 = n C_0 p^0 q^n + n C_1 p^1 q^{n-1} + \dots + n C_n p^n q^0$$

$$\text{mean} = np$$

$$\text{Variance} = npq. \text{ i.e } np(1-p)$$

Bernoulli's distribution: Special case of binomial distribution where only one trial is conducted.

so mean = p

Variance = $p(1-p)$ or pq .

- so for Bernoulli distribution mean is always the probability of success.
- variance is probability of success and probability of failure so variance is always less than mean.

For Bernoulli Distribution PMF

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

for $x=0, 1$ (as success or failure)

$$\begin{aligned} P(X=1) &= p^1 (1-p)^{1-1} \\ &= p \quad = P(\text{success}) \end{aligned}$$

$$\begin{aligned} P(X=0) &= p^0 (1-p)^1 \\ &= 1-p \quad = P(\text{failure}) \end{aligned}$$

Ex:- Approximately 1 in 200 American adults are lawyers. One American adult is randomly selected. What is the distribution of the number of lawyers?

here

- 1 adult selected \rightarrow one trial
- either lawyer or not lawyer - only two results
- probability of one selection being a lawyer is $\frac{1}{200}$ and not dependent on other person.
- All Bernoulli Trial conditions are satisfied.

$$\therefore P(X=x) = \left(\frac{1}{200}\right)^x \times \left(1 - \frac{1}{200}\right)^{1-x}$$

$$x = 0, 1.$$

$$P(X=1) = \frac{1}{200}$$

$$P(X=0) = \frac{199}{200}$$

thus we can logically also but Bernoulli gives nice explanation.

Other distributions built on independent Bernoulli trials are:

- 1) Binomial Distribution - number of successes in independent Bernoulli trials
- 2) Geometric Distribution - Distribution of trials to get the first success in independent Bernoulli trials.
- 3) Negative binomial distribution - It is the distribution of no. of trials to get the r th success in independent Bernoulli trials.

what is Probability Distribution used for ?

- The probability distribution is one of the important concepts in statistics.
- It has huge applications in business, engineering, medicine and other major sectors. It is majorly used to make future predictions based on a sample for a random experiment.
- Example in business sector it is used to predict if there will be profit or loss to the company using any new strategy or by proving any hypothesis test in the medical field, etc.

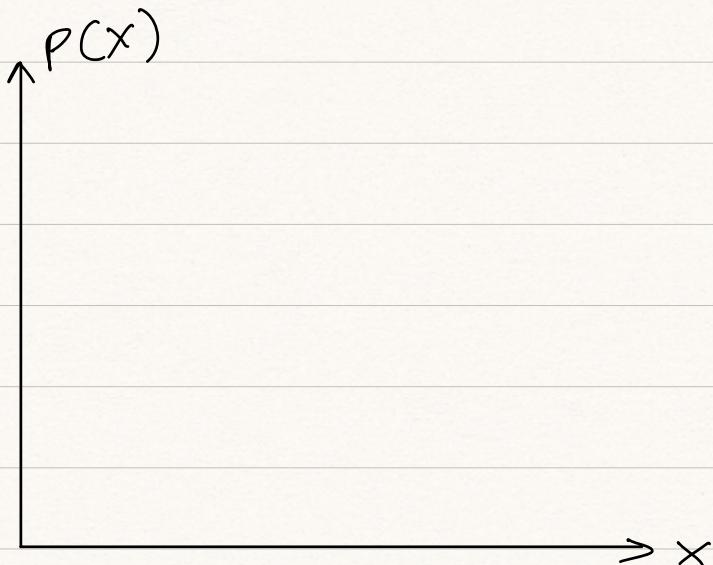
what is the importance of Probability distribution in statistics ?

- In statistics, to estimate the probability of a certain event to occur or estimate the change in occurrence and random phenomena modelled based on the distribution.
- Statisticians take a sample of the population to estimate the probability of occurrence of an event.

Probability distribution of a Discrete Random

Variable:

we can plot this considering the Random Variable on x-axis and probability of the Random Variable on y-axis.



ex:- If we toss 2 coins and assign random variable number of heads

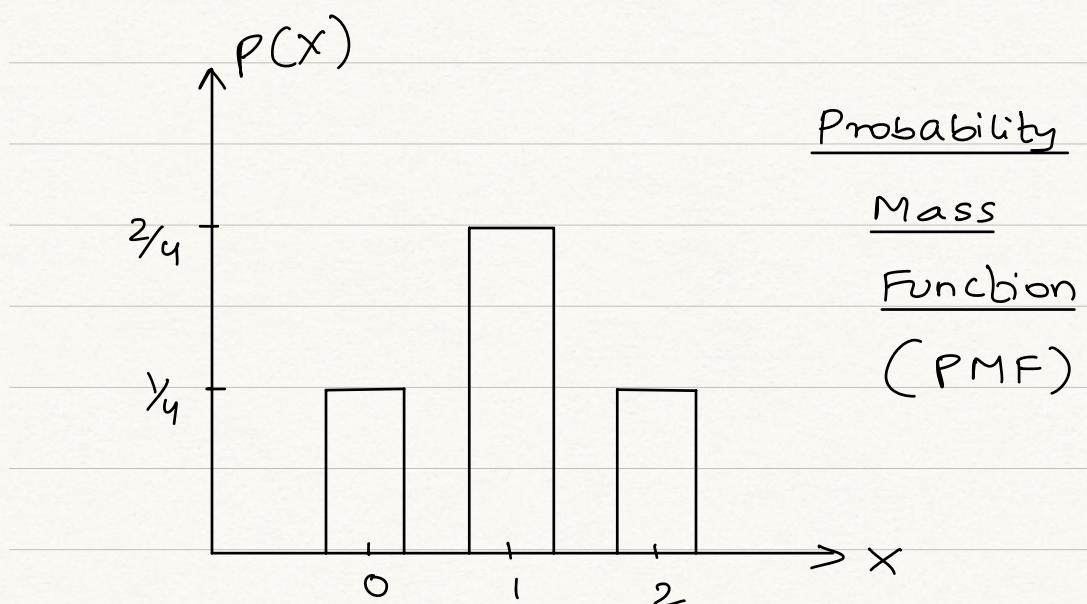
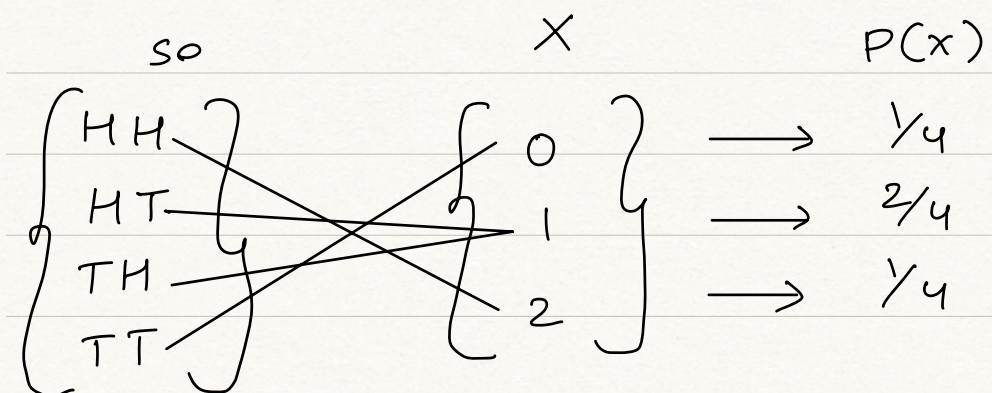
$$S = \{ \begin{matrix} HH \\ HT \\ TH \\ TT \end{matrix} \}$$

$$X - HH \xrightarrow{\text{output}} Z$$

$$x - HT \rightarrow 1$$

X — TH → I

$$x - \pi\pi \longrightarrow 0$$



For Discrete Random Variable when we plot

Random Variable vs Probability of Random

Variable we call it probability Mass function.

It gives the probability that a discrete random

Variable is exactly equal to some value. Sometimes it is also known as discrete density function.

Cumulative Distribution Function: (CDF)

The cumulative distribution function of a real valued random variable X , evaluated at x , is the probability function that X will take a value less than or equal to x .

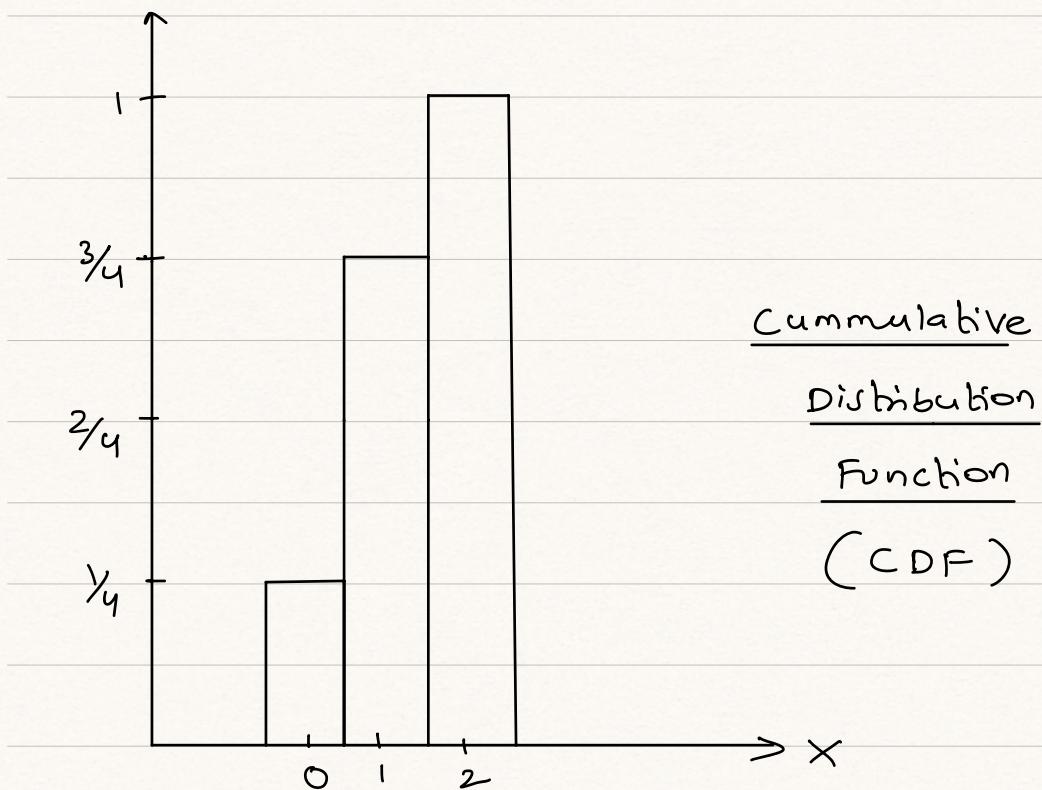
- It is used to describe probability distribution of random variables in a table.
- with the help of these we can plot a CDF plot.
- CDF finds the cumulative probability for given value.
- It can be used for both Discrete Random Variables and Continuous Random Variables.
- It is used to compare the probabilities between values under certain conditions.
- For discrete distribution functions, CDF gives the probability values till what we specify.
- For continuous distribution functions, it gives the area under the probability density function up

to given value specified.

* Cumulative implies to adding previous values to the current one. Similar to calculating the cumulative frequencies.

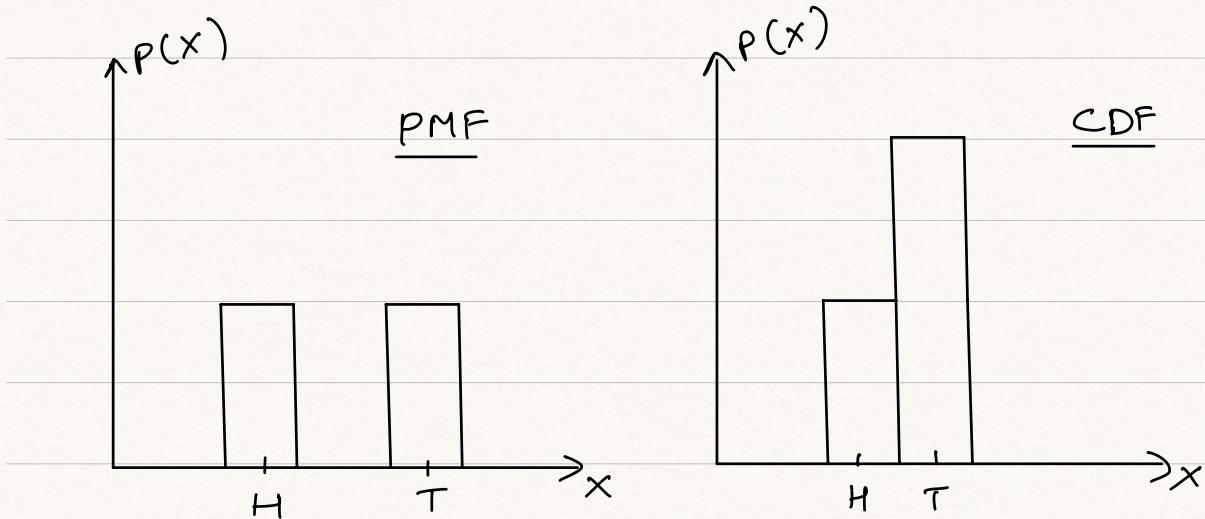
ex:- In toss of two coins

X	0	1	2
$P(X)$	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{4}$
$F(x)$ Cumulative Probability	$\frac{1}{4}$ $= \frac{1}{4}$	$\frac{1}{4} + \frac{2}{4}$ $= \frac{3}{4}$	$\frac{1}{4} + \frac{2}{4} + \frac{1}{4}$ $= 1$

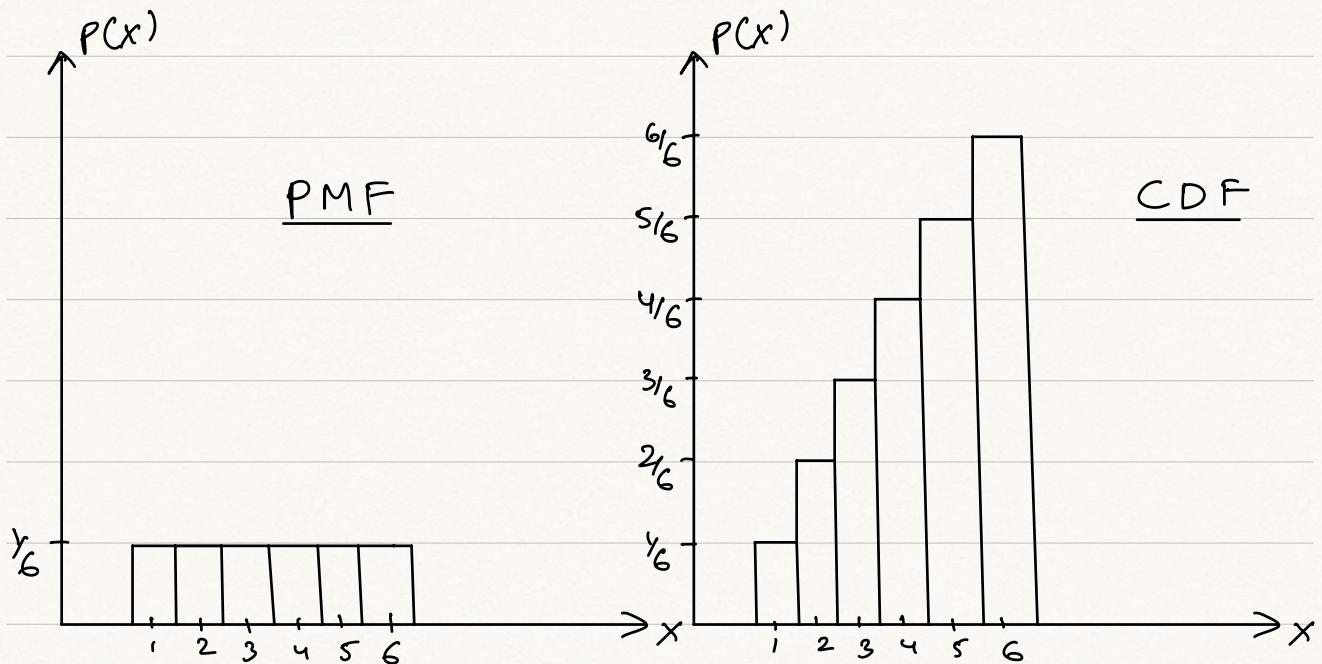


so if we want to find at most or max 1 head i.e $\{HT, TH, TT\} = \frac{3}{4}$
 we can see the same from CDF plot.

ex:- Toss a coin and plot PMF, CDF



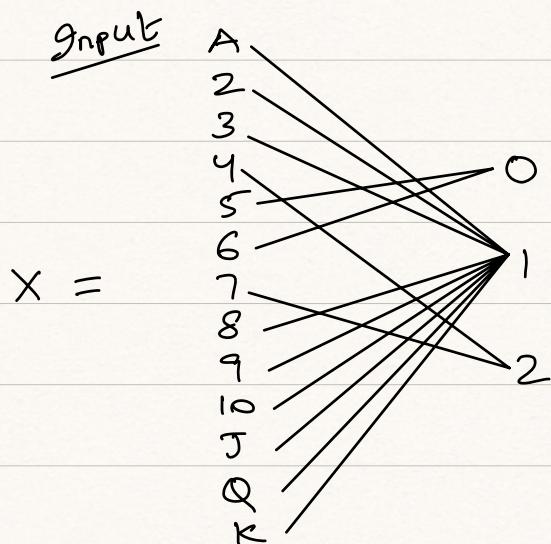
Roll a die



Q) From a deck of spade cards (1s) we have removed 5 & 6 and inserted extra 4 & 7.

Draw PMF and CDF for every event in sample space.

$$\text{Sample space} = \{A, 2, 3, 4, 5, 6, 7, 8, 9, 10, J, Q, K\}$$



$$P(5) = 0$$

$$P(6) = 0$$

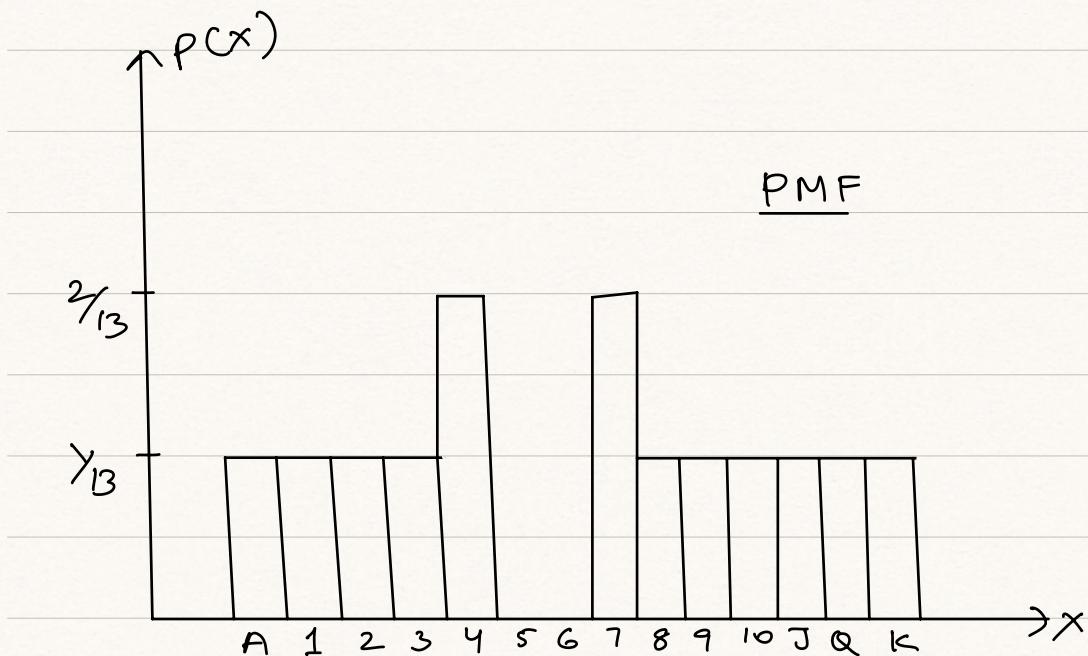
$$P(4) = 2/13$$

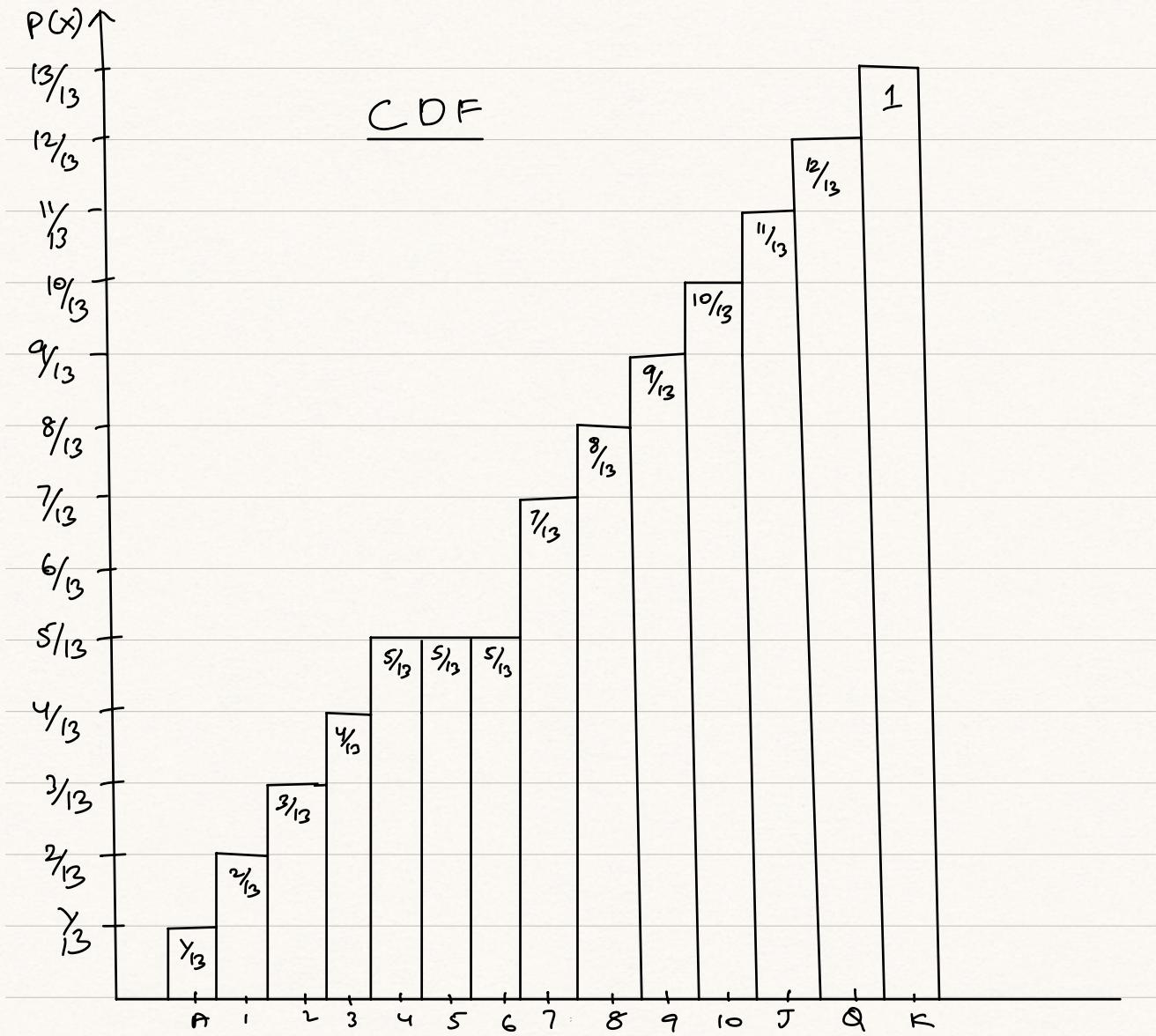
$$P(7) = 2/13$$

for rest of cards

$$P(x) = 1/13$$

$$\text{total cards} = 13$$





Bernoulli's Distribution PMF & CDF

Here the random variable has only two values success (1) or failure (0).

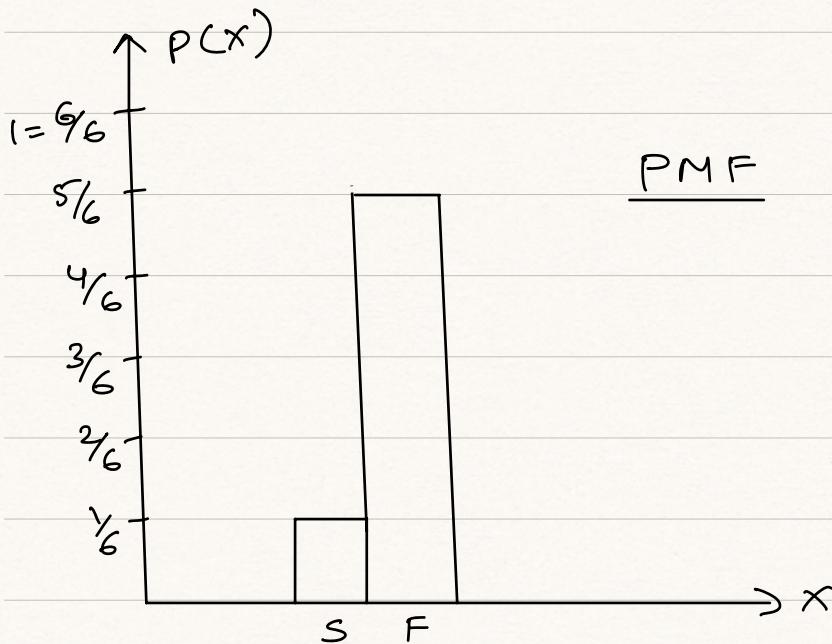
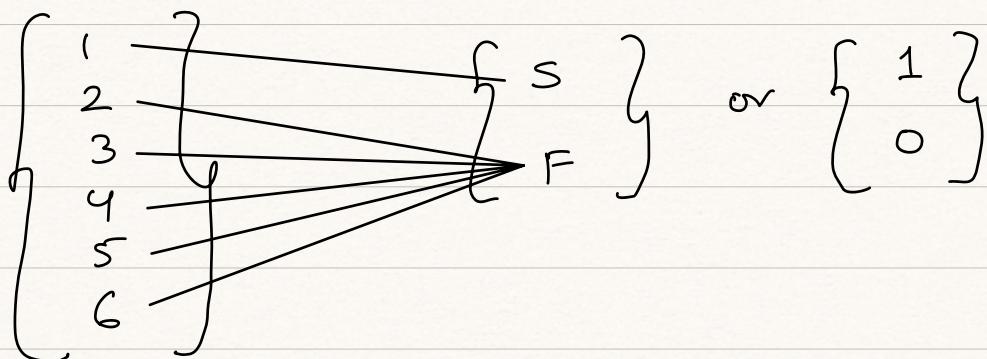
ex:- Getting 1 in toss of a die

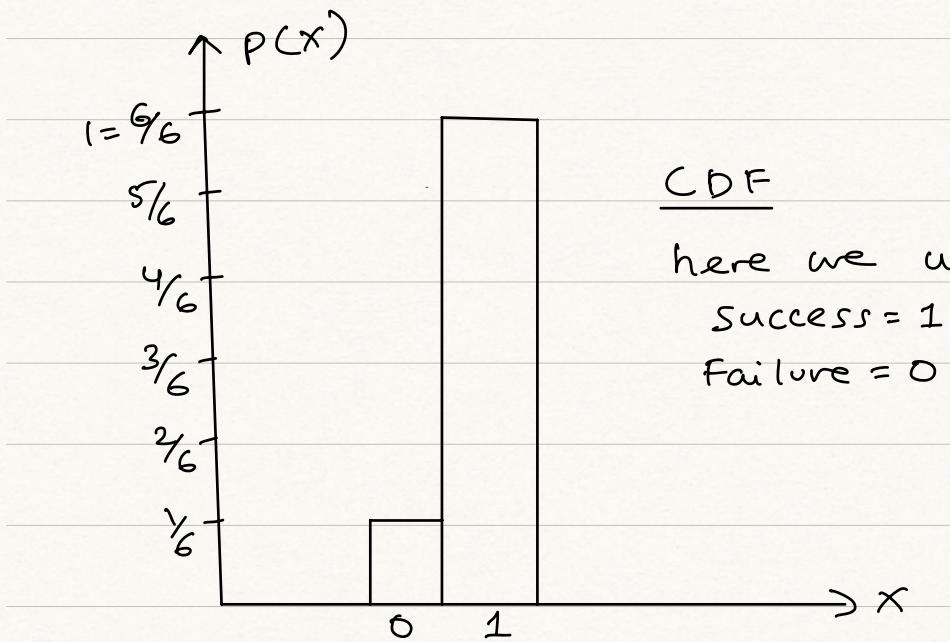
$$P(\text{success}) = 1/6$$

$$P(\text{failure}) = 5/6$$

Sample
space

output

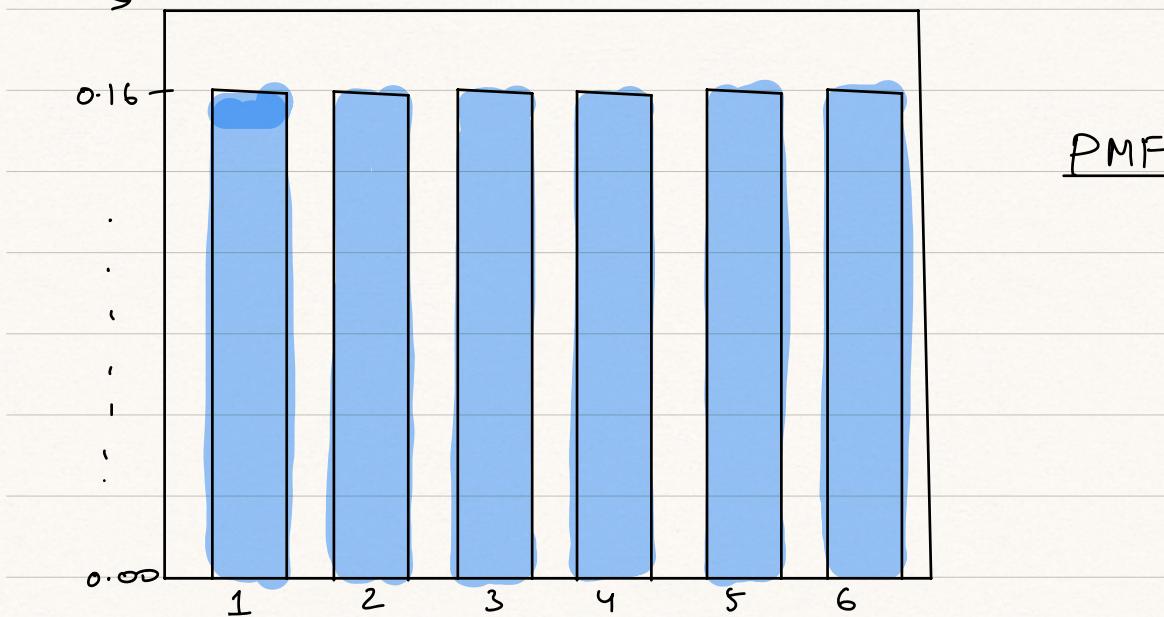




* we can plot this using python.

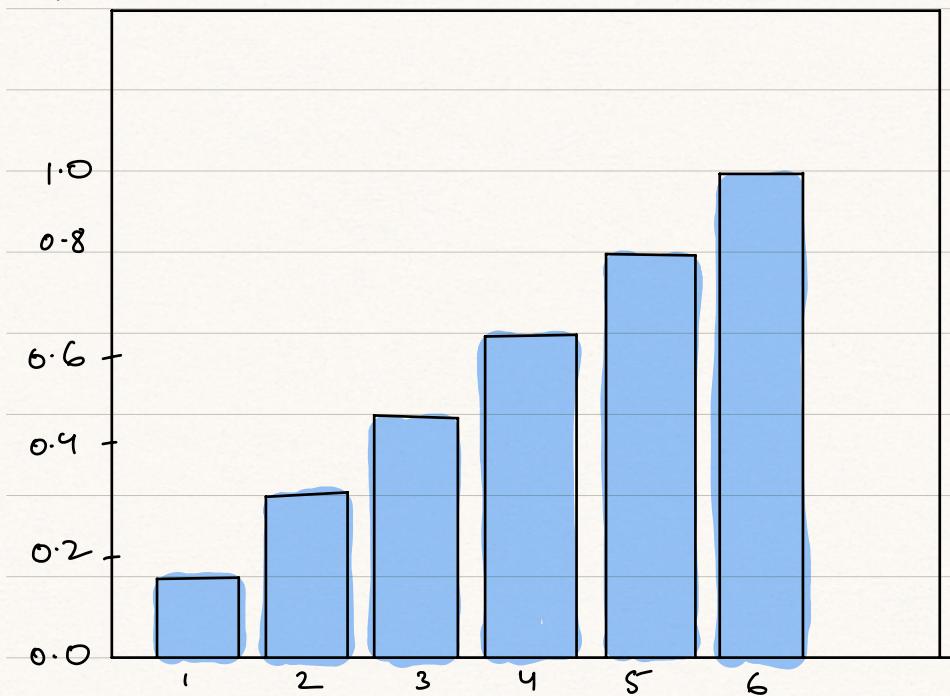
ex:- faces, PFace = np.arange(1,7), np.repeat(1/6, 6)

sns.barplot(x=faces, y=PFace, color="blue")



Cumulative Distribution function

Sns.barplot(x=faces, y=np.cumsum(pFace), color="blue")



If we want to plot PMF & CDF and find mean
variance and standard deviation using python.

ex:-

a = pd.DataFrame({`goals`=[0, 1, 2, 3, 4],
`P(X)`=[0.18, 0.34, 0.35, 0.11, 0.02]})

a

	goals	P(X)
0	0	0.18
1	1	0.34
2	2	0.35
3	3	0.11
4	4	0.02

$$b = a \cdot T$$

b

↳ goals p(x)

0	0	0.18
1	1	0.34
2	2	0.35
3	3	0.11
4	4	0.02

$$\text{mean} = \sum x * p(x)$$

To calculate mean

we use zip function and create zip object and

then create list using the zip object

use for loop and iterate over the list and

find mean using summation of goals * p(x).

$$\text{mean} = 0$$

for i, j in list(zip(b['goals'], b['p(x')]))) :

$$\text{mean} += i * j$$

print(mean)



$$1.4500\cdots 2$$

$$\text{Variance} = \sum (x - \mu)^2 p(x)$$

Here μ is mean calculated above

$$\text{Variance} = 0$$

$$\text{Mean} = 1.4500..2$$

for i,j in list(zip(b['goals'], b['P(x)'])):

$$\text{Variance} += ((1-\text{mean})^{**2}) * j$$

Print(Variance)

$$\downarrow 0.9475$$

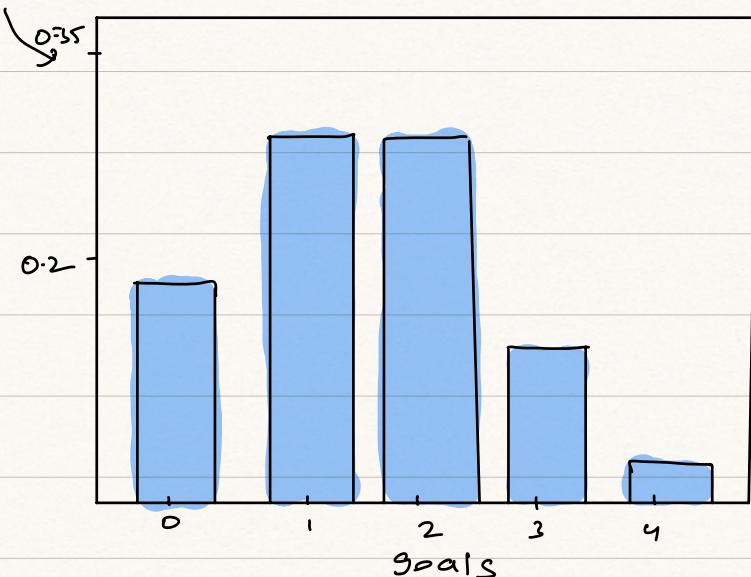
$$\text{std} = \text{abs}(\text{variance}^{**0.5})$$

print(std)

$$\downarrow 0.973396..$$

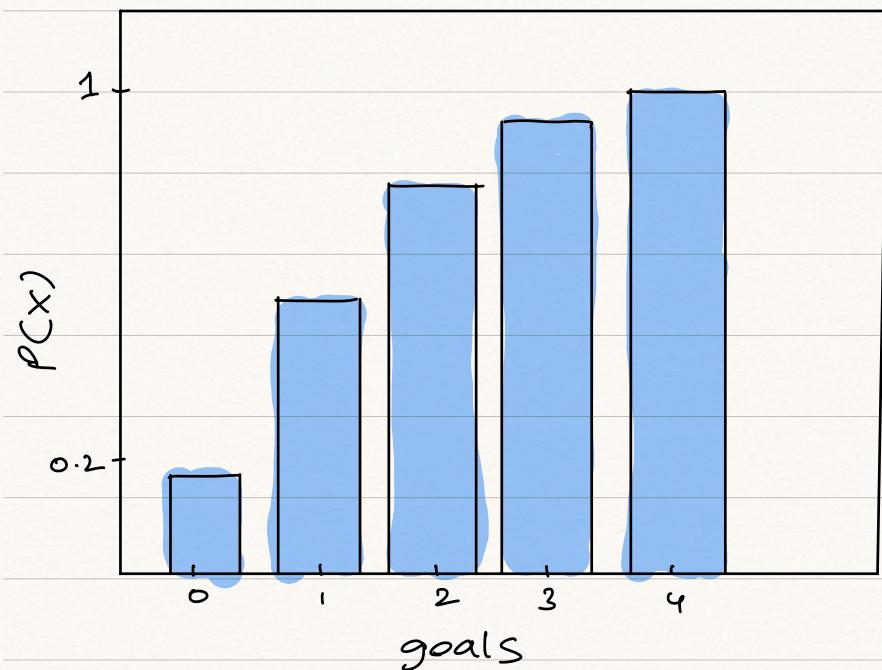
PMF plot

sns.barplot(data=a, x="goals", y="P(x)", color='blue')



CDF plot

Sns.barplot(x=a["goals"], y=np.cumsum(a['P(x)']), color='blue')



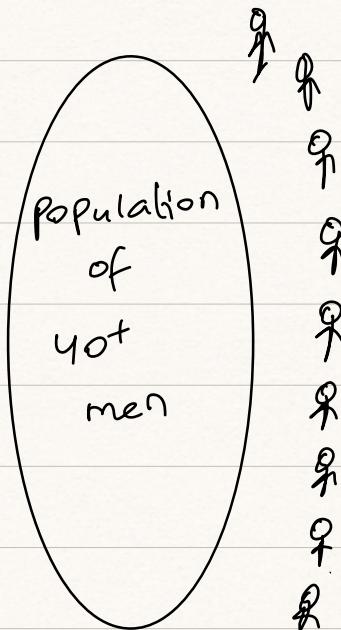
Binomial - refers to population so not finite distribution and with replacement.

Hypergeometric - refers to a finite sample and distribution without replacement

Ex :- Binomial Distribution with coding

Probability of obesity in age 40+ men is 0.014

- (i) if we take a sample of 10 people what is the probability of all of them being obese



Since we are dealing sample of any 10 people from the population. No replacement so independent trials.

$$P(e_1 \cap e_2 \cap e_3 \cap \dots \cap e_{10}) = (P(e))^10$$

as $P(e_1) = P(e_2) = \dots = P(e)$

$$\text{so } P(\text{all 10 being obese}) = (0.014)^{10} \\ = 2.8925 \times 10^{-19}$$

- (ii) if we pick 2 people out of 10 people what is the probability of one being obese and other being not obese

here choose r from n so ${}^n C_r$

$$= {}^{10} C_2 = \frac{10!}{8! 2!} = 45$$

so we can pick 2 people from 10 people in
45 ways.

Now in these one being obese and other
not being obese is

$$45 \times (0.014) ? ?$$

From 10 People:

$$\begin{aligned} P('0' \text{ being obese}) &= {}^{10} C_0 (0.014)^0 \times (1-0.014)^{10-0} \\ &= 1 \times 1 \times (0.986)^{10} \\ &= 0.8684 \end{aligned}$$

$$\begin{aligned} P('10' \text{ being obese}) &= {}^{10} C_{10} \times (0.014)^{10} (1-0.014)^{10-10} \\ &= (0.014)^{10} \end{aligned}$$

$$= 2.8925 \times 10^{-19}$$

$$\begin{aligned} P('1' \text{ being obese}) &= {}^{10} C_1 (0.014)^1 \times (1-0.014)^{10-1} \\ &= 0.1233 \end{aligned}$$

$$P(\text{exactly 3 being obese}) = {}^{10}C_3 (0.014)^3 (0.986)^7$$
$$= 0.00029833$$

$$P(2 \text{ being obese}) = {}^{10}C_2 (0.014)^2 (0.986)^8$$
$$= 0.007879$$

$$P(\min 3 \text{ being obese}) = 1 - (P(0) + P(1) + P(2))$$
$$\chi = 3, 4, 5, \dots, 10$$
$$= 1 - (0.8684 + 0.1233 + 0.007879)$$
$$= 0.00042077$$

$$P(7 \text{ men being obese}) = 1 - P(\min 3 \text{ being obese})$$
$$\chi = 7, 8, 9, 10$$
$$= 1 - 0.00042077$$
$$= 0.99957923$$

$$P(\text{exactly 7 being obese}) = {}^{10}C_7 (0.014)^7 \times (0.986)^3$$
$$= 0.12126 \times e^{-10}$$

Coding of Binomial distribution using Python classes:-

```
import math
```

```
from math import factorial as fact
```

```
class Binomial:
```

```
    def __init__(self, population_size, p_s):
```

```
        self.n = population_size
```

```
        self.p = p_s
```

```
        self.dtable = self.binomialDistTable()
```

```
        self.dft = self.dtable.T
```

```
    def binomialDistTable(self):
```

```
        l = []
```

```
        for x in range(0, self.n + 1):
```

```
            choice = fact(self.n) / ((fact(self.n - x)) * fact(x))
```

```
            ps = self.p ** x
```

```
            pf = (1 - self.p) ** (self.n - x)
```

```
            px = choice * ps * pf
```

```
            l.append((x, px))
```

```
        return pd.DataFrame(l, columns=['x', 'prob-x'])
```

```
    def min_of(self, num):
```

```
        if num > 0 and num <= self.n:
```

```
return self.dtable.iloc[num:, 1].values.  
cumsum()[-1]
```

else:

```
print(f'num {num} does not qualify')
```

```
def max_of(self, num):
```

```
if num > 0 and num <= self.n:
```

```
return self.dtable.iloc[:num+1, 1].  
values.cumsum()[-1]
```

else:

```
print(f'num {num} does not qualify')
```

```
def pmfplot(self):
```

```
return sns.barplot(data=self.dtable, x='x',  
y='prob-x')
```

```
def cdfplot(self):
```

```
return sns.barplot(data=self.dtable, x='x',  
y=self.dtable['prob-x'].values.cumsum())
```

```
# Creating Binomial class objective
```

```
Bd = Binomial(10, 0.014)
```

```
# Probability of min 3 being obese
```

```
Bd.min_of(3)
```

→ 0.00030587

Bd. max-of (3)

→ 0.999992459

Bd. dtable

	X	Prob-X
0	0	$8.68.. \times e^{-01}$
1	1	--
2	2	--
3	3	$2.9833.. \times e^0$
4	4	--
5	5	--
6	6	--
7	7	--
8	8	--
9	9	--
10	10	--

Probability of exactly getting $x=3$

`BD::dtable::loc[3, "Prob-x"]`

$\rightarrow 0.00029833$

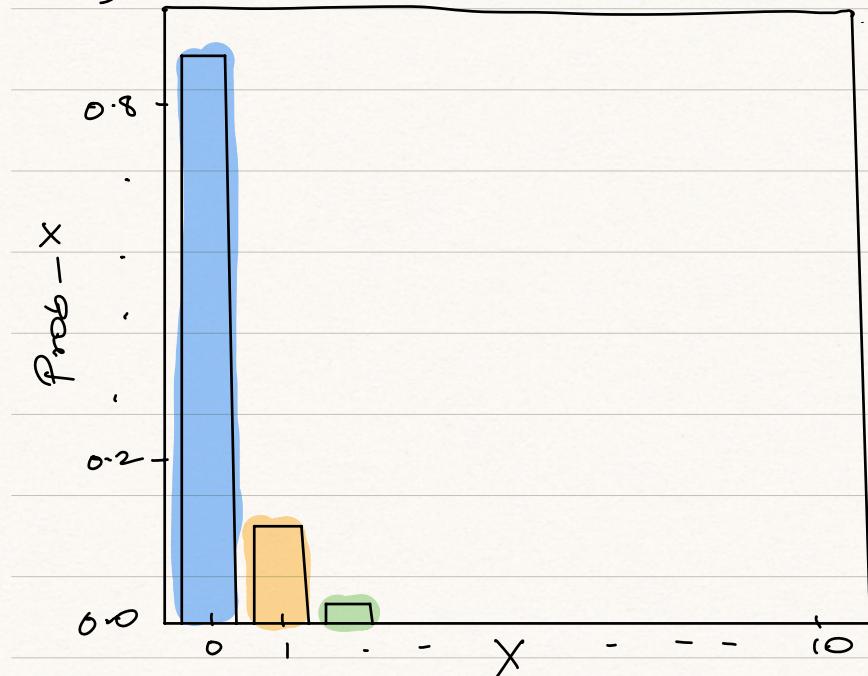
Probability Distribution table

Bd.dft

	0	1	2	3	4	5	6	7	8	9	10
X	0.00..	1.00..	.	-	.	.	-	-	-	-	-
prob-X	0.86..	0.12..	-	.	-	.	-	-	-	-	.

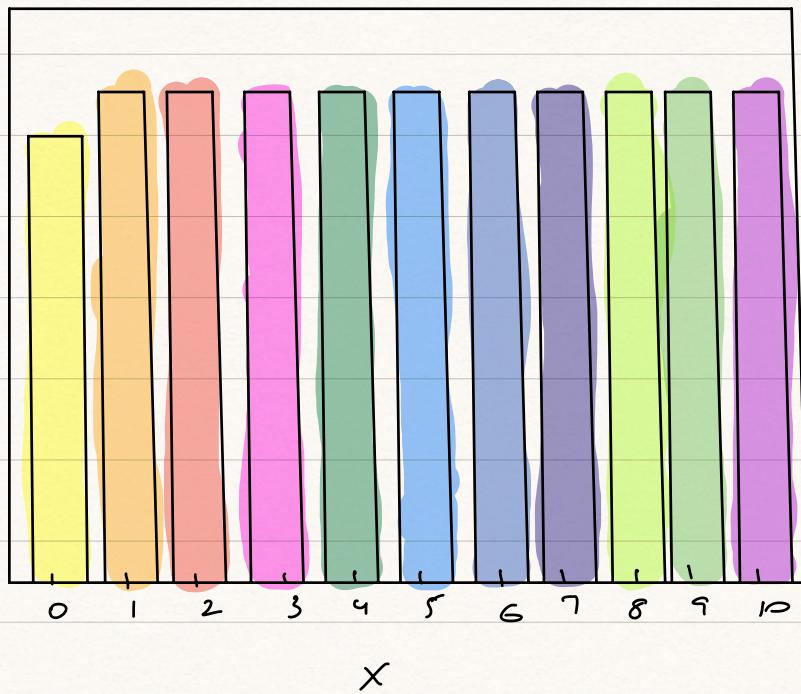
Pmf plot for Binomial

Bd.pmfPlot()



cdf plot for Binomial

Bd.cdfPlot()



we have a scipy stats function called

Scipy.stats.binom.pmf(x,n,p)

or

Scipy.stats.binom.cdf(x,n,p)

Hypergeometric distribution.

- Another Discrete probability Distribution

- we should already know

$$n_C_r = \frac{n!}{(n-r)! r!}$$

- Hypergeometric distribution is not based on Bernoulli trials as without replacement so no longer independent trials.

- Main difference from Binomial distribution is this is done without replacement.

Let us consider an example first for better understanding.

Ex:- An urn contains 6 red balls and 14 yellow balls. 5 balls are randomly drawn without replacement. What is the probability exactly 4 red balls are drawn.

So once the ball is drawn then it is placed aside. So the probability of the next draw depends on the first one so no longer independent trials so Binomial cannot be used.

Let us calculate using logic before deriving formulas.

If we are randomly selecting 5 balls then any sample of 5 balls is equally likely, so:

$$P(\text{Exactly 4 red balls})$$

$$= \frac{\# \text{ of samples that result in 4 red balls and 1 yellow ball}}{\# \text{ possible samples of size 5}}$$

$$= \frac{6 \times {}^4C_4 \times {}^{14}C_1}{{}^{20}C_5} \quad \begin{array}{l} \text{from 6 red balls pick 4 red balls} \\ \text{from 14 yellow balls pick 1 yellow (not red) ball.} \end{array}$$

$$= 0.01354$$

here if we think logically

$$\text{on first trial } P(\text{Red}) = 6/20$$

$$\text{on second trial if first ball is Red } P(\text{Red}) = 5/19$$

so $P(\text{success})$ keeps changing so not independent.

For us to use Binomial distribution the question should be

ex:- An urn contains 6 red balls and 14 yellow balls. 5 balls are randomly drawn with replacement. what is the probability exactly 4 red balls are drawn.

$P(\text{success}) = \frac{6}{20}$ and it will not change

$$P(4 \text{ red balls}) = {}^5C_4 \left(\frac{6}{20}\right)^4 \times \left(\frac{14}{20}\right)^1 \\ = 0.02835.$$

different from Hypergeometric probability.

Hypergeometric distribution

- Suppose we are randomly sampling n objects without replacement from a source of N objects that contains a successes and $N-a$ failures
- X represents the number of successes in the sample.

Then X has a hypergeometric distribution:

(PMF)

$$P(X=x) = \frac{a_{Cx} * (N-a)_{C(n-x)}}{N_{Cn}}$$

or
 $P(x)$

Here we see
 $N-a+a = N$
 $x+n-x = n$
we can use for
double check

for $x = \text{Max}(0, n-(N-a)), \dots, \min(a, n)$

i.e first we need to pick x from total successes of a i.e a_{Cx}

then we need to pick $n-x$ failures from total $N-a$ failures

Denominator is total samples i.e picking n objects from N objects.

Here x is # of successes

it cannot be bigger than # of objects sampled n and also cannot be bigger than the total successes i.e $\min(a, n)$

The min value x can take cannot be lesser than 0 or n minus total failures i.e $n - (N-a)$
so min value of x is $\max(0, n - (N-a))$

mean of Hypergeometric distribution

$$M = n \times \frac{a}{N}$$

\uparrow no. of samples \uparrow proportion of successes in the total population N

$$\text{Variance} = \sigma^2 = V(X) = \left(\frac{N-n}{N-1}\right) * n * \left(\frac{a}{N}\right) \left(\frac{N-a}{N}\right)$$

ex:- suppose a large high school has 1100 female students and 900 male students.

A random sample of 10 students is drawn.
what is the probability exactly 7 of the students selected are female?

Ans) Though not mentioned this is without replacement

$$P(X=7) = ??$$

logically:- denominator: total samples
so picking 10 from 2000
i.e $2000_{C_{10}}$.

Numerator: probability of picking exactly 7 female students is from 1100 we pick 7 i.e 1100_{C_7} and remaining 3 students we pick from 900 male students

$$\text{i.e } 900_{C_3} \text{ i.e } 1100_{C_7} \times 900_{C_3}$$

$$\therefore P(X=7) = \frac{1100_{C_7} \times 900_{C_3}}{2000_{C_{10}}} \quad \begin{array}{l} 1100+900=2000 \\ 7+3=10 \\ \text{double check satisfied} \end{array}$$

$$= 0.166490$$

if formula $N=2000$

$$n=10$$

$$a=1100$$

$$x=7$$

$$P(X=7) = \frac{1100C_7 \times (2000-1100)C_{10-7}}{2000C_{10}}$$

$$= 0.166490.$$

If we assume this is with replacement and use Binomial distribution instead

$$\text{so } P(\text{success}) = \frac{1100}{2000}$$

$$P(X=7) = 10C_7 \left(\frac{1100}{2000}\right)^7 \times \left(1 - \frac{1100}{2000}\right)^{10-7}$$

$$= 0.166478$$

we can see that this is pretty close to probability using Hypogeometric distribution.

* The binomial distribution can sometimes provide a reasonable approximation to the hypogeometric distribution.

Rough guideline: If we are not sampling more than say 5% of the population, the binomial provides a reasonable approximation.

So we use binomial distribution where ever we can approximate for easier calculations.

Multi-variate hypergeometric

Example: Suppose a business employs 12 democrats, 24 republicans, and 8 independents.
i.e $12 + 24 + 8 = 44$ total.

If a random sample of 6 employees is drawn, what is the probability there are 3 democrats, 2 republicans and 1 independent in the sample.

$$P(x) = \frac{12C_3 \times 24C_2 \times 8C_1}{44C_6}$$

so we can apply hypergeometric methods to more than two types of object and this is sometimes called multi-variate Hypergeometric distribution.

Hypogeometric Distribution with coding

Simplified formula used for coding is

$$P(X=x) = \frac{n! k! (N-n)! (N-k)!}{N! x! (k-x)! (n-x)! (N-k-n+x)!}$$

here

N - Total size of sample/population

k - total size of interest group

n - total number picked

x - input of random variable.

Hypogeometric distribution using python classes:-

```
import math
```

```
from math import factorial as fact
```

```
class HypoGeometric:
```

```
    def __init__(self, population_size, p_s):
```

```
        self.n = population_size
```

```
        self.p = p_s
```

```
        self.dtable = self.hypoGeometricDistTable()
```

```
        self.dft = self.dtable.T
```

```

def hypergeometricDistTable(self):
    l = []
    for x in range(0, self.k+1):
        p_x = (fact(self.n) * fact(self.k) *
                fact(self.N - self.n) * fact(self.N -
                self.k)) / (fact(self.n) * fact(x) *
                fact(self.k - x) * fact(self.n - x) *
                fact(self.N - self.k - self.n + x))
        l.append((x, p_x))
    return pd.DataFrame(l, columns=['x', 'prob-x'])

```

```
def min_of(self, num):  
    if num > 0 and num <= self.n:  
        return self.dtable.iloc[num:, 1].values.  
        cumsum()[-1]  
  
    else:  
        print(f'num {num} does not qualify')
```

```
def max_of(self, num):  
    if num > 0 and num <= self.n:  
        return self.dtable.iloc[:num+1, 1].values.cumsum()[-1]
```

else:

print(f'num {num} does not qualify')

def pmfplot(self):

return sns.barplot(data=self.dtable, x='x',
y='prob-x')

def cdfplot(self):

return sns.barplot(data=self.dtable, x='x',
y=self.dtable['prob-x'].values.cumsum())

we create a class object

hh = Hypergeometric(230, 5, 10)

probability of min 3 being obese

.min_of(3)

→ 0.0005721...

.max_of(3)

→ 0.999990...

hg.dtable

	x	prob-x
0	0	$7.99 \dots \times e^{-01}$
1	1	--
2	2	--
3	3	$5.6307 \dots \times e^{-04}$
4	4	--
5	5	--

Probability of exactly getting $x=3$

.dtable.loc[3, 'Prob-x']

$$\hookrightarrow 5.6307 \times e^{-04}$$

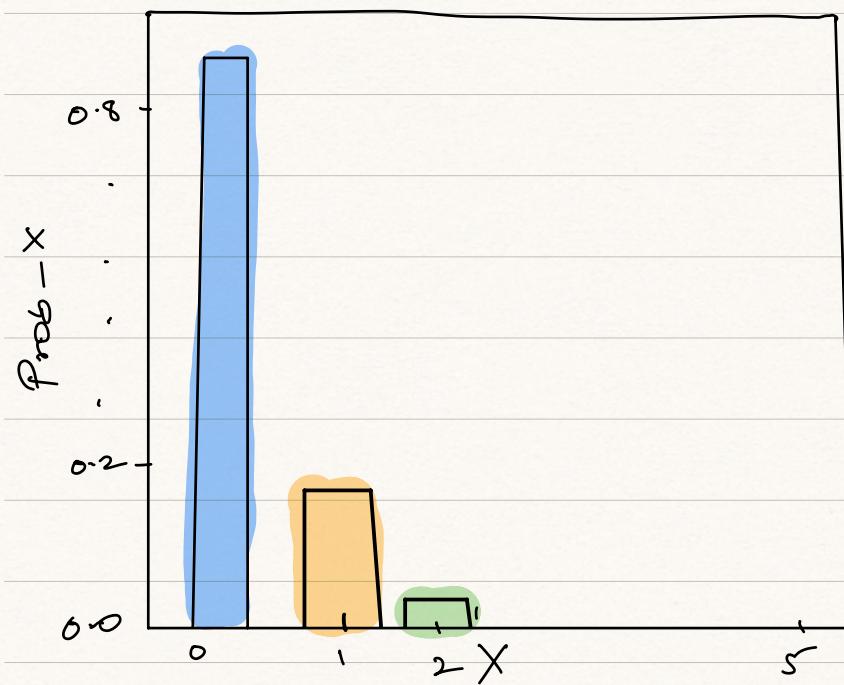
Probability Distribution table

hg.dft	0	1	2	3	4	5
x	0.00..	1.00..	.	-	-	.
prob-x	0.799..	0.184..	-	.	-	.

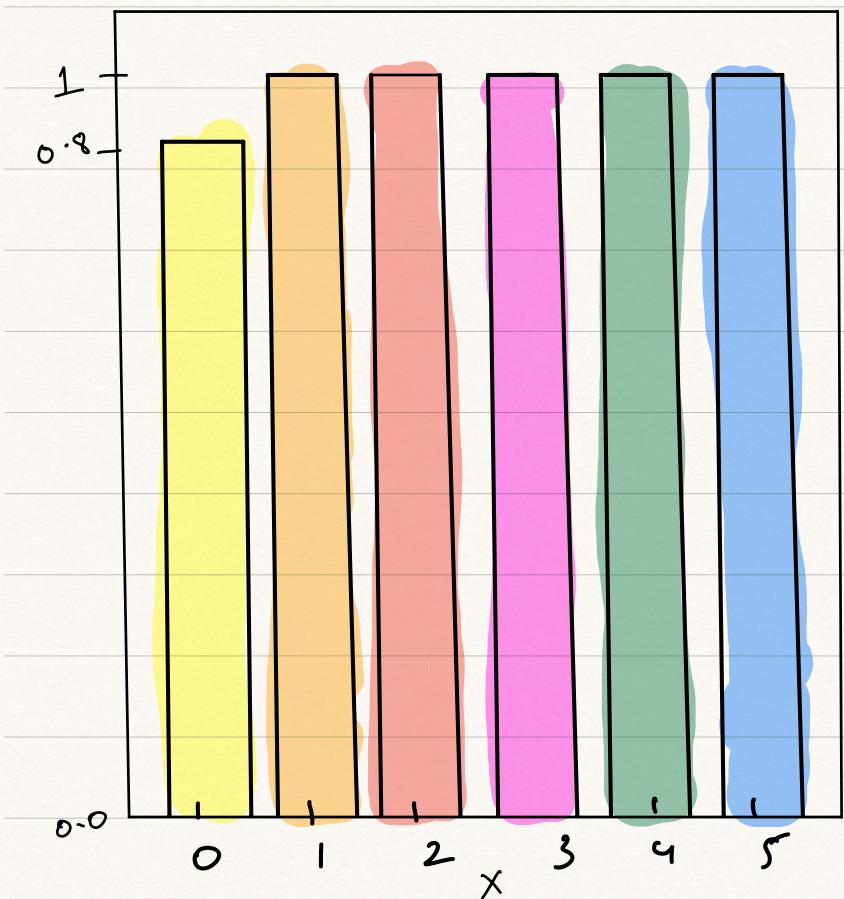
Pmf plot for hypergeometric

.pmfPlot()





#cdf plot for hypogeometric



In scipy we have a stats hypergeometric function called

scipy.stats.hypergeom()

logical explanation of pmf function of hypergeometric function.

PMF of $= \frac{\{a\text{ successes}\} \times \{n-a\text{ failures}\}}{\text{choices of } n \text{ out of } N}$
hyper geometric function

$= \frac{\text{no. of ways of choosing a group of } n \text{ with } a \text{ successes}}{\text{no. of ways to choose } n \text{ from } N}$.

Geometric distribution: To better understand let us consider an example:

- In a large population of adults, 30% have received CPR training.

If adults from this population are randomly selected, what is the probability that the 6th person sampled is the first that have received CPR training?

We can use geometric distribution to calculate this probability.

The geometric distribution is the distribution of the number of trials needed to get the first success in independent repeated Bernoulli trials.

Repeated independent Bernoulli's trials:

SUPPOSE:

- There are independent trials
- Each trial results in one of two possible mutually exclusive outcomes, labelled success and failure.

On any individual trial:

- $P(\text{success}) = p$ and this stays constant from trial to trial

- $P(\text{Failure}) = 1-p$

- X represents the number of trials needed to get the first success.

pmf

for the first success to occur on x^{th} trial:

1. The first $x-1$ trials must be failures $\rightarrow (1-p)^{x-1}$
2. The x^{th} trial must be success $\rightarrow p$

Since trials are independent.

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

$x = 1, 2, 3, \dots$ (no upper bound)

$$\text{mean} = \mu = \frac{1}{p}$$

$$\text{Variance} = \sigma^2 = \frac{1-p}{p^2}$$

so in geometric distribution we just need probability of success.

so if we go back to the example.

- In a large population of adults, 30% have received CPR training.

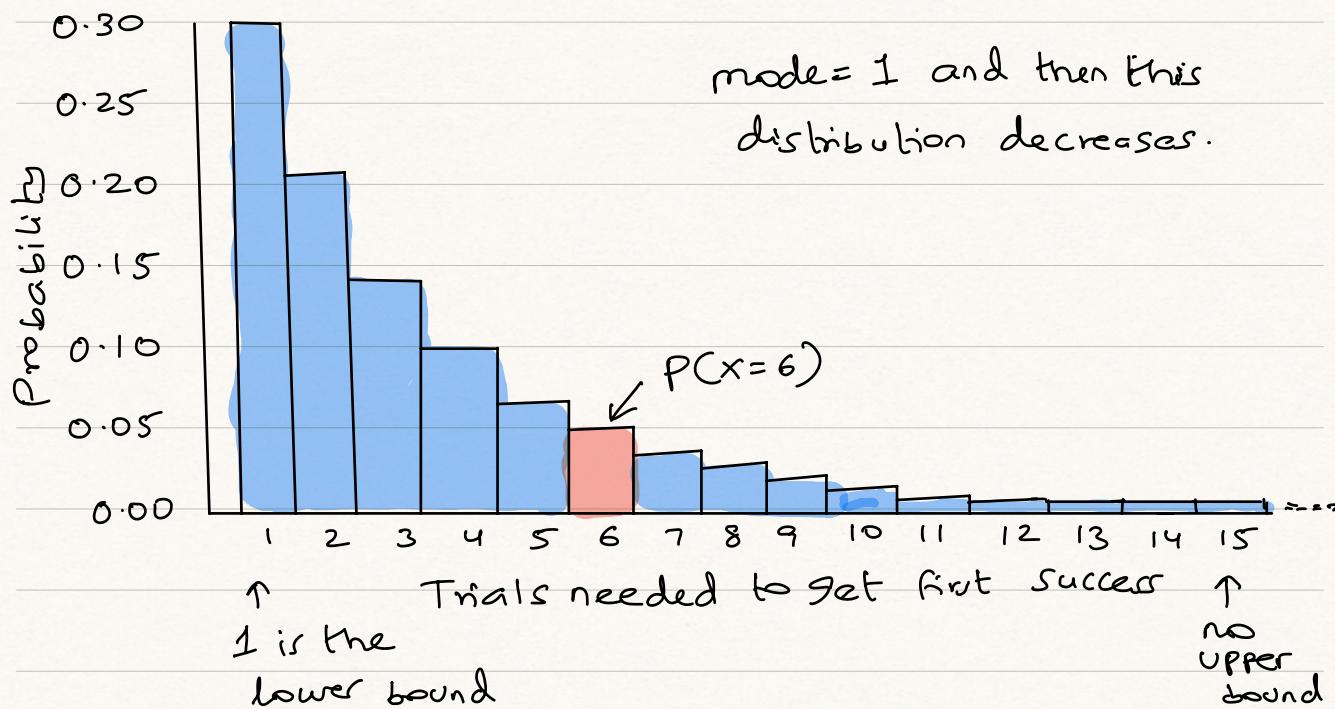
If adults from this population are randomly selected, what is the probability that the 6th person sampled is the first that have received CPR training?

$$p = 30\% = 0.3 \quad \text{on any individual trial}$$

$$n = 6.$$

$$\begin{aligned} P(X=6) &= \text{first 5 not CPR} \times 6^{\text{th}} \text{ CPR training.} \\ &= (1-0.3)^{6-1} (0.3) \\ &= 0.0504 \end{aligned}$$

if we plot this



Geometric distribution always has right skewed distribution.

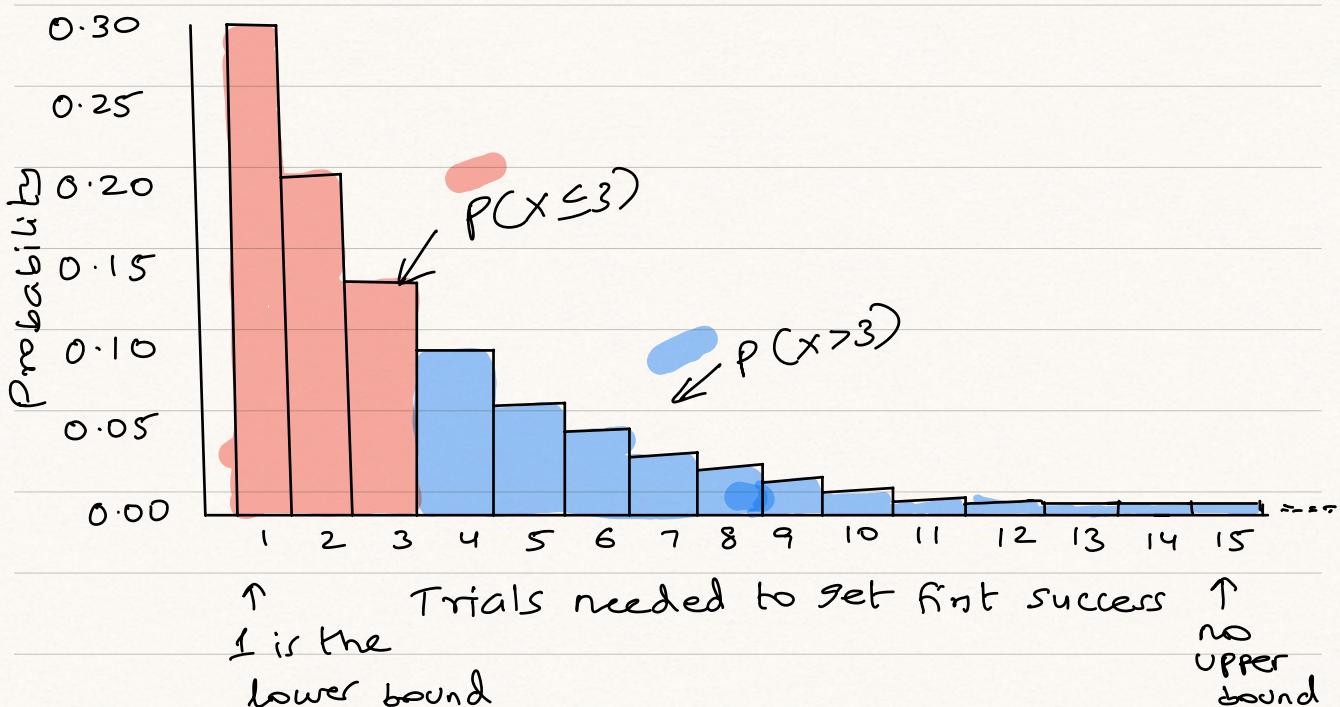
$$\text{mean} = \frac{1}{0.3} = 3.3$$

$$\text{Variance} = \frac{1 - 0.3}{0.3^2} = 7.7$$

$$\sigma = \sqrt{7.7} = 2.79$$

Q) what is the probability that the first person trained in CPR occurs on or before the 3rd person sampled?

$$P(X \leq 3) = ?$$



$$\begin{aligned}
 P(X \leq 3) &= P(X=1) + P(X=2) + P(X=3) \\
 &= 0.7^0 \times 0.3 + 0.7^1 \times 0.3 + 0.7^2 \times 0.3 \\
 &= 0.657
 \end{aligned}$$

or (we use when there are large no. of probabilities)

$$\begin{aligned}
 P(X \leq 3) &= 1 - P(X > 3) \\
 &\hookrightarrow \text{i.e. getting three failures.} \\
 &= 1 - 0.7^3 \\
 &= 0.65
 \end{aligned}$$

The cumulative distribution function for the geometric distribution:

$$F(x) = P(X \leq x) = 1 - (1-p)^x$$

for $x = 1, 2, 3, \dots$

Poisson distribution

Suppose we are counting the number of occurrences of an event in a given unit of time, distance, area or volume.

For example:

- The number of car accidents in a day.
- The number of dandelions in a square metre plot of land

The number of events is going to be a random variable that may or may not have the Poisson distribution, depending on the specifics of the situation.

* But a Poisson distribution is a count of the number of occurrences of an event.

If we consider in terms of time

Suppose:

- Events are occurring independently i.e knowing when one event happens has absolutely no information on when another event will occur.
- And the probability that an event occurs in a given length of time does not change

through time i.e the theoretical rate at which the events are occurring does not change through time.

so we can say the events are occurring randomly and independently.

* If these conditions hold then the random variable x which represents the no. of events in a fixed unit of time has Poisson distribution.

* Till now we were looking at random Variables that take finite values. In Poisson distribution the random variable can take values up to ∞ . Though the probability near infinity would be minute.

* Here the mean and variance are same and represented as λ

$$\mu = \sigma^2 = \lambda$$

PMF of Poisson distribution that we will use to calculate probabilities is

$$P(x=x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

$$n = 0, 1, 2, \dots, \infty$$

e is base of natural logarithm

$$e \approx 2.71828$$

it is an irrational number that has infinite decimal places

ex:- plutonium - 239 is an isotope of plutonium that is used in nuclear weapons and reactors.

One nanogram of plutonium - 239 will have an average of 2.3 radioactive decays per second, and the number of decays will follow a poisson distribution.

what is the probability that in a 2 second period there are exactly 3 radioactive decays?

let x represent the number of decays in a 2 second period.

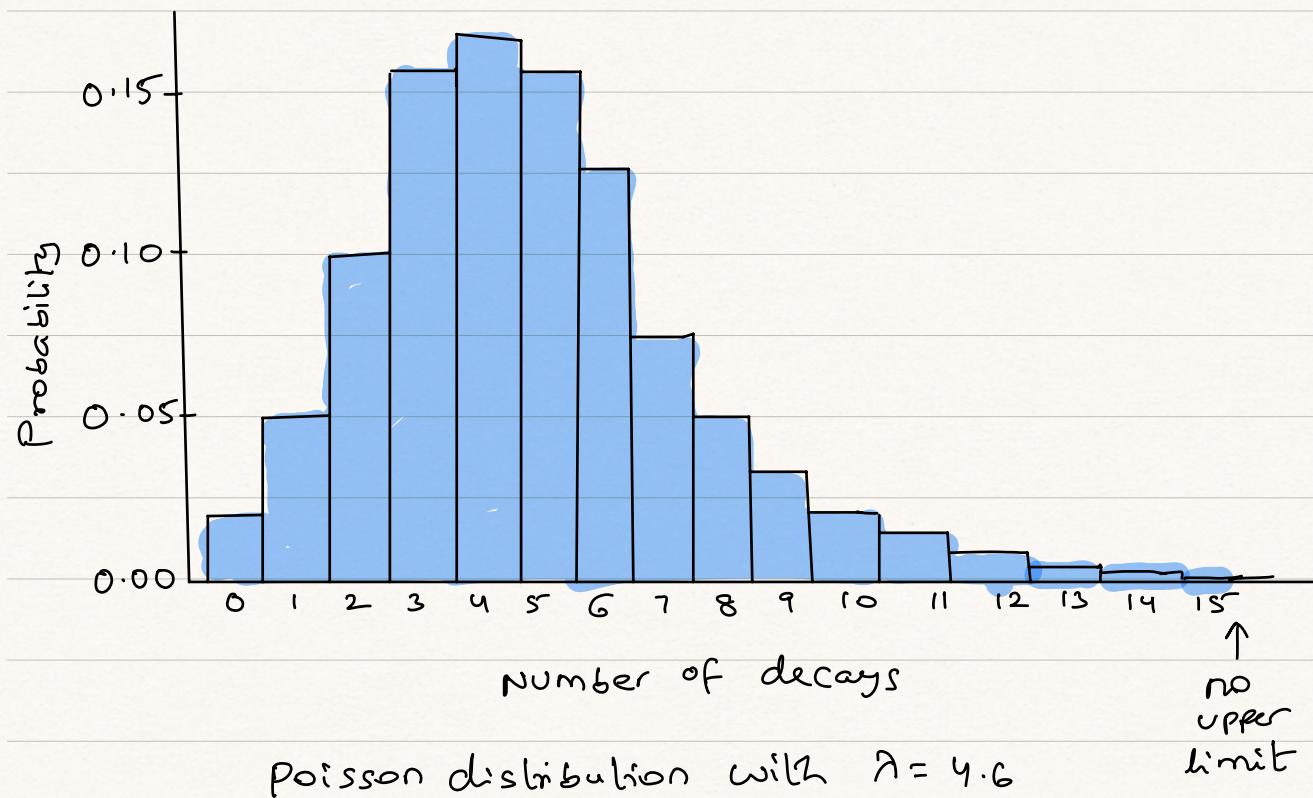
$$\text{mean for 2 second} = 2.3 \times 2$$

$$\lambda = 4.6$$

$$P(X=3) = \frac{4.6^3 e^{-4.6}}{3!}$$

$$= 0.163$$

if we were to calculate probabilities for the different possible values x can take and plot them we will get the probability distribution plot



$$\mu = \lambda = 4.6$$

$$\sigma^2 = \lambda = 4.6$$

$$\sigma = \sqrt{4.6}$$

* when λ is large the distribution is close to symmetric.

* when λ is close to 0, the right skewness can

be pretty strong

$$P(X \leq 3) = ?$$

$$= P(X=0) + P(X=1) + P(X=2) + P(X=3)$$

$$= \frac{4.6^0 e^{-4.6}}{0!} + \frac{4.6^1 e^{-4.6}}{1!} + \frac{4.6^2 e^{-4.6}}{2!} + \frac{4.6^3 e^{-4.6}}{3!}$$

$$= 0.010 + 0.046 + 0.106 + 0.163$$

$$= 0.326$$

The relationship between Binomial and Poisson

Distributions

- The binomial distribution tends toward the Poisson distribution as $n \rightarrow \infty$, $p \rightarrow 0$ and np stays constant.
- The poisson distribution with $\lambda = np$ closely approximates the binomial distribution if n is large and p is small.
 - this is why the radioactivity decays of plutonium has a poisson distribution. Even for a tiny bit of plutonium, there are a very large number of atoms, and each one has a tiny probability of experiencing a radioactive decay in a two second period. So in the example we just worked through, it was in its underlying nature a binomial problem with a very large n and very small p . And that's why the number of radioactive decays is very well approximated by poisson distribution.
- The poisson distribution is typically used as an approximation to the true underlying reality.
- It can be difficult to determine whether a

random variable has a Poisson distribution.