

PRACTICAL FILE OF PROBABILITY FOR COMPUTING

RAMANUJAN COLLEGE



UNIVERSITY OF DELHI

DSC 06: PROBABILITY FOR COMPUTING

2024-25

SEMESTER 2

SUBMITTED BY

NAME:

GURPREET SINGH PAUL

ROLL NO.

24570022(24020570023)

COURSE:

B.Sc(H)Computer Science

SEMESTER:2

SUBMITTED TO

Dr Aakash

Assistance Professor

(Operational Research),

Department of Computer

Science, Ramanujan College,

University of Delhi, CR Park

Main Road, Block H, Kalkaji,

New Delhi-110019

INDEX

S.No.	TOPIC	Pg.No.
	Plotting and fitting of Binomial distribution and graphical representation of probabilities.	
	Plotting and fitting of Multinomial distribution and graphical representation of probabilities.	
	Plotting and fitting of Poisson distribution and graphical representation of probabilities.	
	Plotting and fitting of Geometric distribution and graphical representation of probabilities	
	Plotting and fitting of Uniform distribution and graphical representation of probabilities	
	Plotting and fitting of Exponential distribution and graphical representation of probabilities.	
	Plotting and fitting of Normal distribution and graphical representation of probabilities.	
	Calculation of cumulative distribution functions for Exponential and Normal distribution.	
	Given data from two distributions, find the distance between the distributions.	
	Application problems based on the Binomial distribution	
	Application problems based on the Poisson distribution	
	Application problems based on the Normal distribution	
	Presentation of bivariate data through scatter-plot diagrams and calculations of covariance	
	Calculation of Karl Pearson's correlation coefficients	
	To find the correlation coefficient for a bivariate frequency distribution	
	Generating Random numbers from discrete (Bernoulli, Binomial, Poisson) distributions.	
	Generating Random numbers from continuous (Uniform, Normal) distributions	
	Find the entropy from the given data set.	

Acknowledgement

I would like to express my heartfelt gratitude to everyone who contributed to the successful completion of this practical assignment.

First and foremost, I extend my sincere appreciation to my teacher, **Dr Aakash, Assistance Professor (Operational Research), Ramanujan College, University of Delhi**, whose guidance, patience, and insightful feedback were invaluable throughout this process. Their dedication to fostering a deep understanding of mathematics has greatly enhanced my learning experience.

I am also grateful to my classmates for their encouragement and collaboration, which made the learning journey both engaging and enjoyable.

Name: Gurpreet Singh Paul

Roll no. 24020570023

1. Plotting fitting of Binomial distribution and graphical representation of probabilities.

Binomial Distribution:

The binomial distribution is a discrete probability distribution. It describes the outcome of binary scenarios, e.g. toss of a coin.

Binomial Distribution Formula:

$$P(x;n,p) = {}^n C_x p^x (1-p)^{n-x}$$

Where:

- $P(x; n, p)$ is the probability of x successes in n trials in an experiment which can result in exactly two outcomes (success or failure).
- p is the probability of success on an individual trial.
- n is the number of trials. x is the
- total number of successes.

Mean($\mu=np$)

Variance($\sigma^2=npq$)

Implementation in Excel:

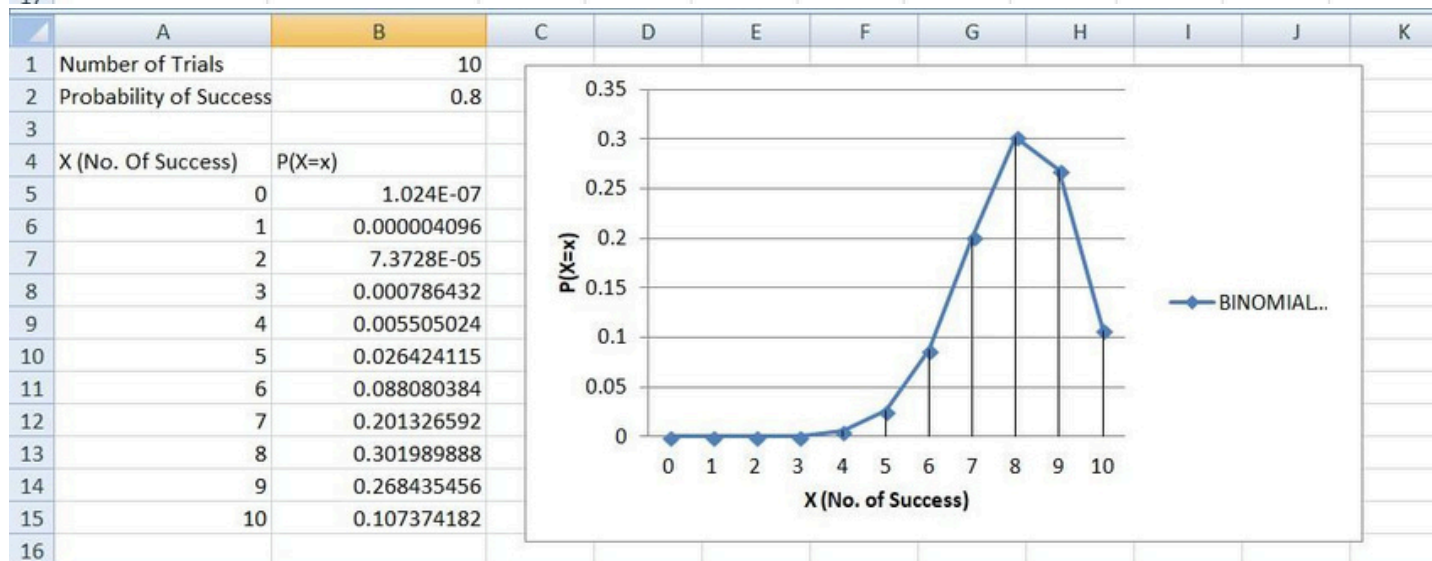
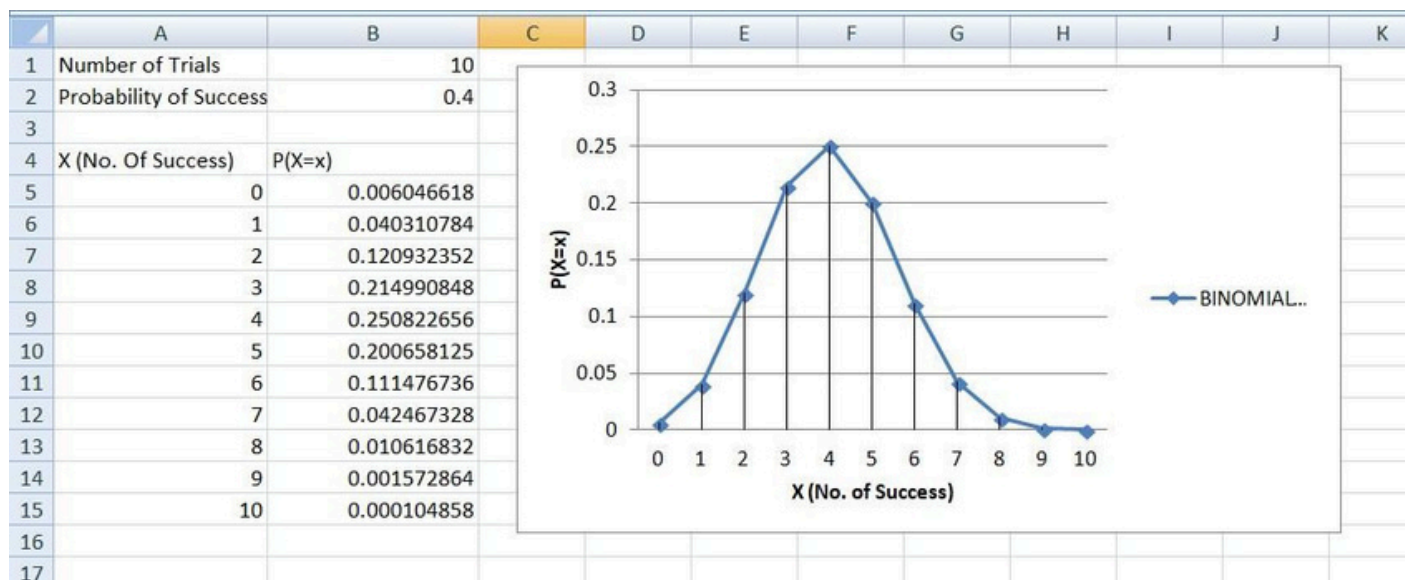
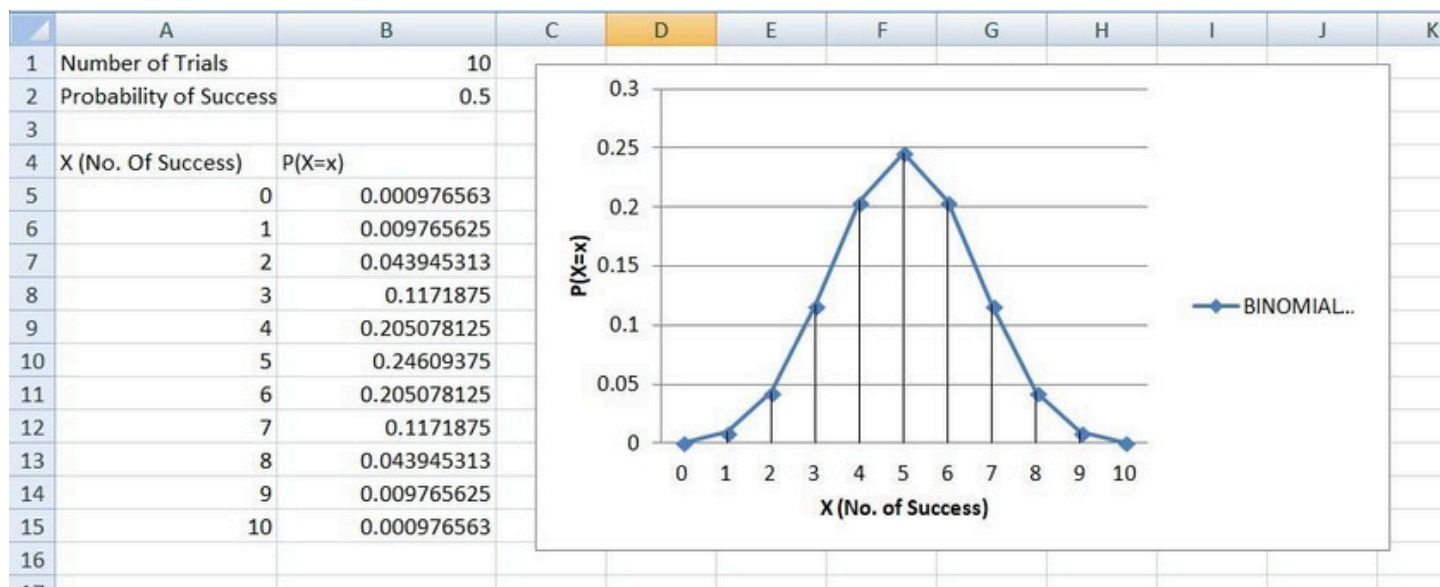
=BINOM.DIST(number_s, trials, probability_s, cumulative)

Where:

- number_s: number of successes. trials:
- total number of trials.
- probability_s: probability of success on each trial. cumulative: TRUE returns the cumulative
- probability; FALSE returns the exact probability

Skewness In Case Binomial Distribution Can Be Defined As Follows:

- If $p=0.5$, the binomial distribution will be symmetrical, regardless of the value of n .
- If $p \neq 0.5$, the distribution will be skewed.
- If $p < 0.5$, the distribution will be positively skewed or right-skewed. This means the bulk of the probability falls in the smaller numbers and the distribution tails off to the right.
- If $p > 0.5$, the distribution will be negatively skewed or left-skewed. This means the bulk of the probability falls in the larger numbers and the distribution tails off to the left.



2. Plotting and fitting of Multinomial distribution and graphical representation of probabilities.

The multinomial distribution is a multivariate generalization of the binomial distribution. Consider a trial that results in exactly one of some fixed finite number k of possible outcomes, with probabilities, p_1, p_2, \dots, p_k (so that $p_i \geq 0$ for $i = 1, \dots, k$ and $\sum_{i=1}^k p_i = 1$), and there are n independent trials. Then let the random variables X_i indicate the number of times outcome number i was observed over the n trials. Then $X = (X_1, X_2, \dots, X_k)$ follows

a multinomial distribution with parameters n and p , where $p = (p_1, p_2, \dots, p_k)$,

Multinomial Distribution Formula

$$p(x_1, x_2, \dots, x_k) = \left[\frac{n!}{x_1! \cdot x_2! \cdot \dots \cdot x_k!} \right] \cdot p_1^{x_1} \cdot p_2^{x_2} \cdot \dots \cdot p_k^{x_k}$$

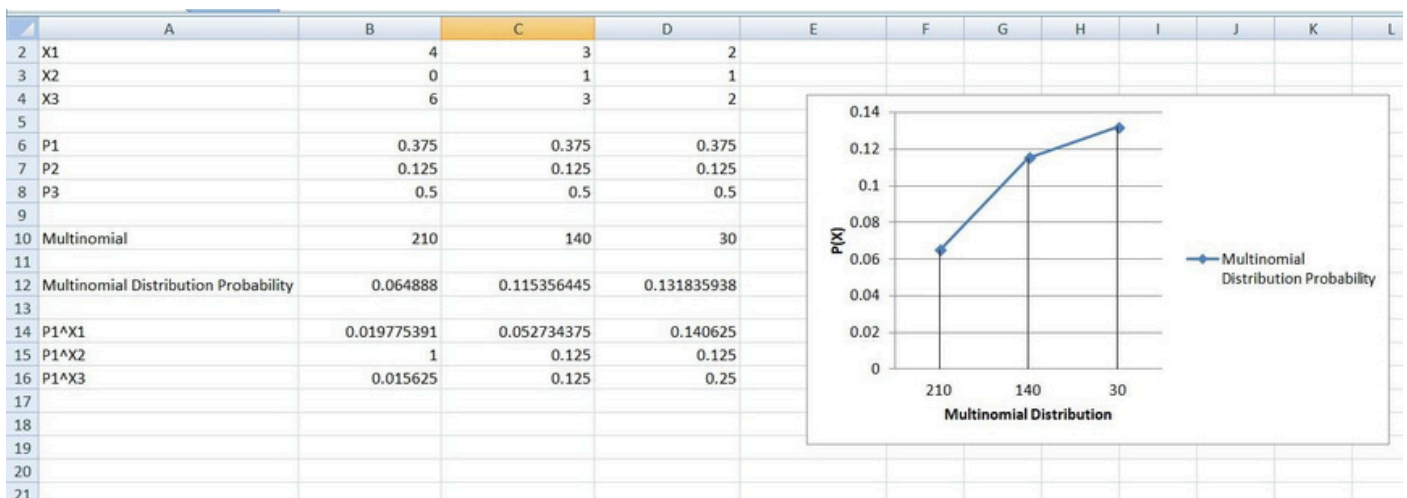
$$\text{Cov}(X_i, X_j) = -np_i p_j \quad (i \neq j)$$

When $X = (x_1, x_2, \dots, x_k)$ follows a multinomial distribution with the PMF given above, $X_{\{i\}}$ follows a binomial distribution with n trials and success probability p_i .

how to implement in excel

Multinomial MULTINOMIAL(X1,X2,X3)

Probability = MULTINOMIAL PRODUCT($p_1^{X_1}, p_2^{X_2}, p_3^{X_3}$)



3. Plotting and fitting of Poisson distribution and graphical representation of probabilities.

The Poisson distribution is a type of discrete probability distribution that determines the likelihood of an event occurring a specific number of times (k) within a designated time or space interval. This distribution is characterized by a single parameter, λ (lambda), representing the average number of occurrences of the event.

Poisson Distribution Formula

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

Mean	$\mu = E(X) = \lambda$
Variance	$\sigma^2 = V(X) = \lambda$
Standard Deviation	$\sigma = \sqrt{\sigma^2} = \sqrt{\lambda}$

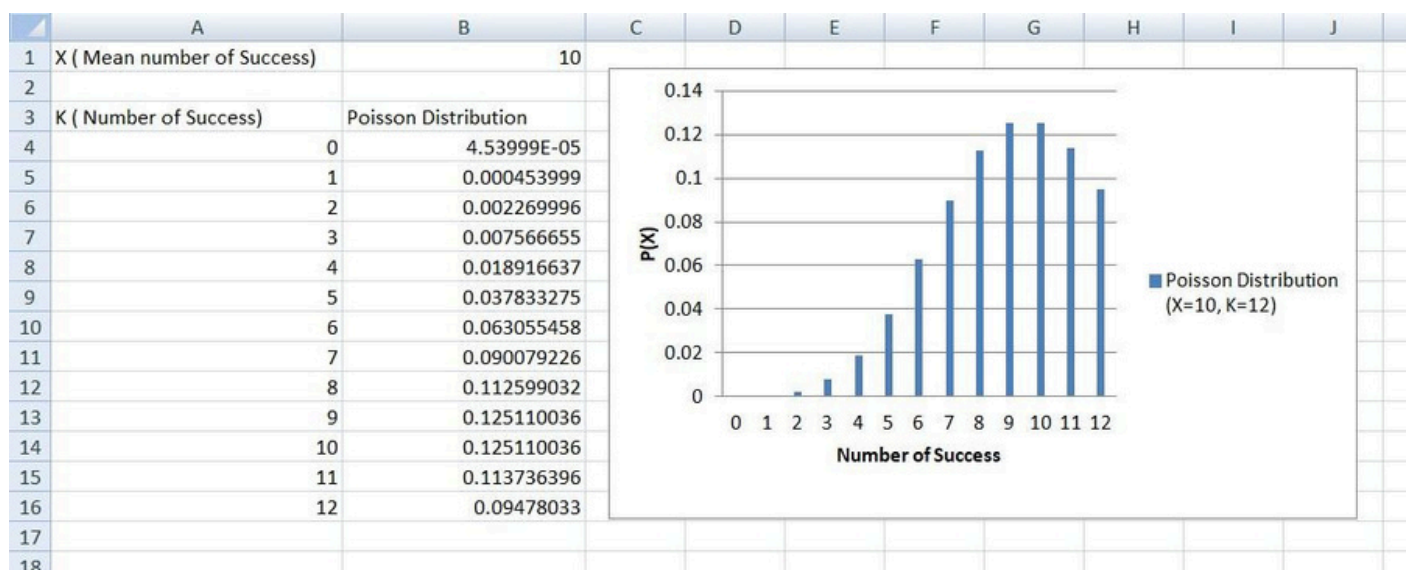
Where:

- $P(X=k)$ is the probability of observing k events
- e is the base of the natural logarithm (approximately 2.71828)
- A mean number of success that occur during a specific interval, λ and k
-

is the number of success how to implement in excel

POISSON.DIST(number_s,average,cumulative)

POISSON.DIST(k, λ , FALSE)



4. Plotting and fitting of Geometric distribution and graphical representation of probabilities.

In a Bernoulli trial, the likelihood of the number of successive failures before success is obtained is represented by a geometric distribution, which is a sort of discrete probability distribution. A Bernoulli trial is a test that can only have one of two outcomes: success or failure. In other words, a Bernoulli trial is repeated until success is obtained and then stopped in geometric distribution.

A geometric distribution is a discrete probability distribution that indicates the likelihood of achieving one's first success after a series of failures. The number of attempts in a geometric distribution can go on indefinitely until the first success is achieved. Geometric distributions are probability distributions that are based on three key assumptions.

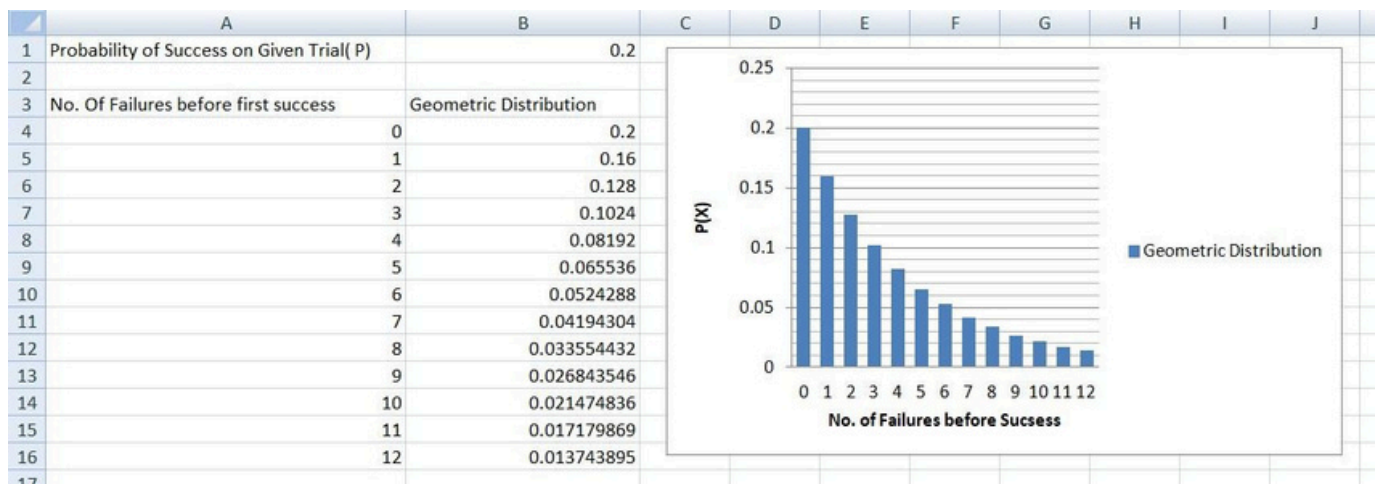
- The trials that are being undertaken are self-contained.
- Each trial may only have one of two outcomes: success or failure.
- For each trial, the success probability, represented by p , is the same

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

Mean:	$\mu = E(X) = \frac{1}{p}$	Geometric Distribution
Variance:	$\sigma^2 = V(X) = \frac{(1-p)}{p^2}$	

Where :

- k : number of failures before first success p :
- probability of success on each trial
- The chance of a trial's success is denoted by p , whereas the likelihood of failure is denoted by q , $q = 1-p$ in this case. $X \sim G(p)$ represents a discrete random variable, X , with a geometric probability distribution.



how to implement in excel
 Probability = $(1-p)^k \cdot p$

5. Plotting and fitting of Uniform distribution and graphical representation of probabilities.

A uniform distribution is a distribution that has constant probability due to equally likely occurring events. It is also known as rectangular distribution (continuous uniform distribution). It has two parameters a and b: a = minimum and maximum. The distribution is written as U(a, b) b =

A uniform distribution is a type of probability distribution where every possible outcome has an equal probability of occurring. This means that all values within a given range are equally likely to be observed.

Uniform Distribution Formula

The probability density function (PDF) of a continuous uniform distribution defines the probability of a random variable falling within a particular interval. For a continuous uniform distribution over the interval [a,b].

$$f(x) = \frac{1}{b-a} \text{ for } a \leq x \leq b$$

$$\text{Mean } \mu = \frac{a+b}{2}$$

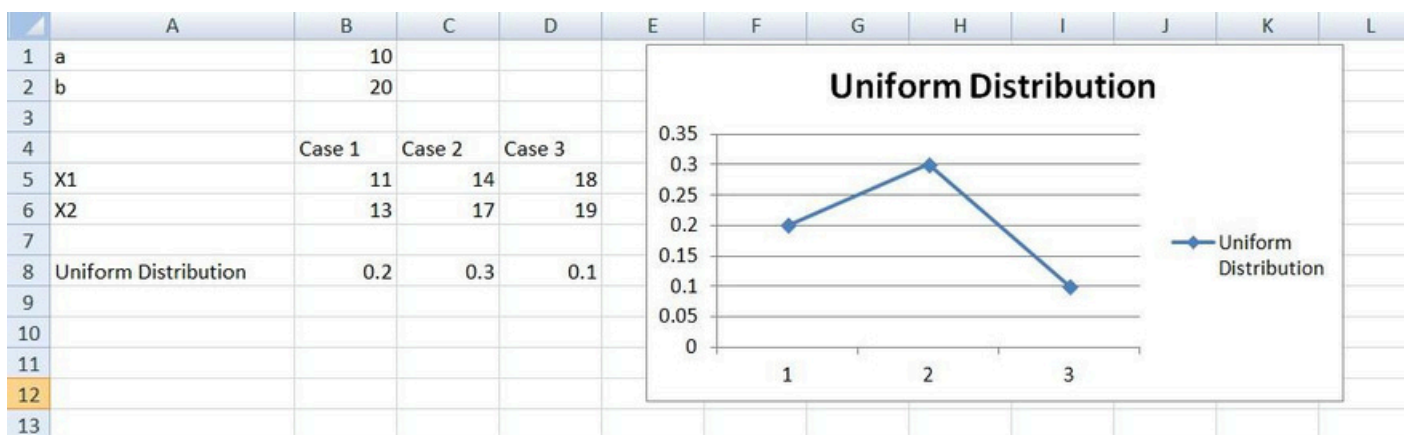
$$\text{Variance } \sigma^2 = \frac{(b-a)^2}{12}$$

how to implement in excel

$$P = (x_2 - x_1) / (b - a)$$

For calculating probability, we need:

1. a: minimum value in the distribution
2. b: maximum value in the distribution
3. x1: the minimum value you're interested in
4. x2: the maximum value you're interested in



6. Plotting and fitting of Exponential distribution and graphical representation of probabilities.

The support (set of values the Random Variable can take) of an Exponential Random Variable is the set of all positive real numbers. Suppose we are posed with the question- How much time do we need to wait before a given event occurs? The answer to this question can be given in probabilistic terms if we model the given problem using the Exponential Distribution, Since the time we need to wait is unknown, we can think of it as a Random Variable, If the probability of the event happening in a given interval is proportional to the length of the interval, then the Random Variable has an exponential distribution. The support (set of values the Random Variable can take) of an Exponential Random Variable is the set of all positive real numbers.

This distribution can be used to solve following type of real life problems-

- How long does a shop owner need to wait until a customer enter a shop.
- How long will a battery continue to work before it dies.
- How long will a computer continue to work before it breakdown.

$$f(x) = \begin{cases} 0, & x < 0 \\ \lambda e^{-\lambda x}, & x \geq 0 \end{cases}$$

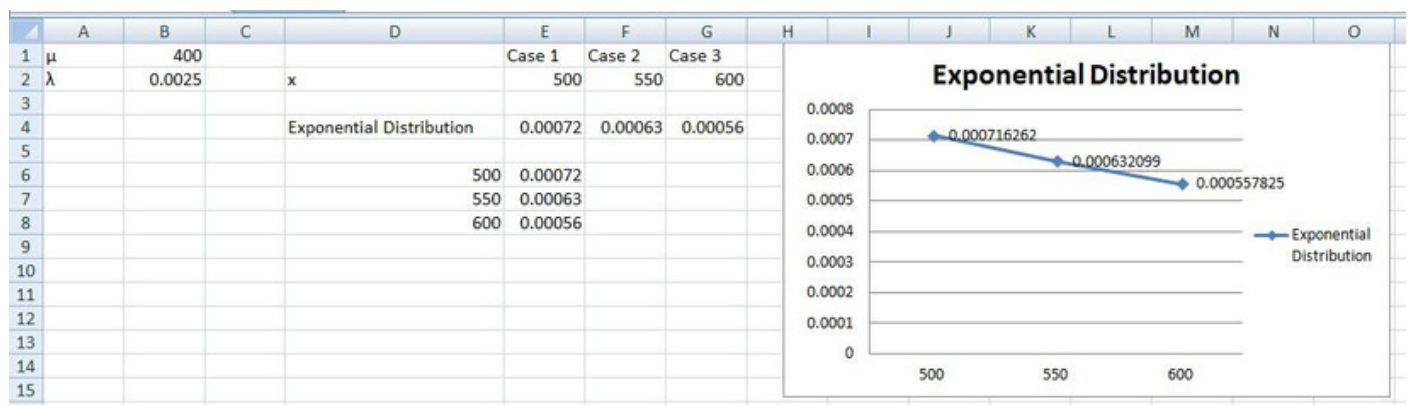
$$E(X) = \frac{1}{\lambda}, \quad \text{Var}(X) = \frac{1}{\lambda^2}$$

Here is the rate parameter and its effects on the density function. e is a constant roughly equal to 2.718

How to Implement in excel

EXPON.DIST(X,lambda,cumulative)

EXPON.DIST (X,lambda, FALSE)



7. Plotting and fitting of Normal distribution and graphical representation of probabilities.

We define Normal Distribution as the probability density function of any continuous random variable for any given system. Now for defining Normal Distribution suppose we take $f(x)$ as the probability density function for any random variable X .

$$f(x) \geq 0 \quad \forall \quad x \in (-\infty, +\infty),$$

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

where, x is Random

- Variable μ is Mean σ
- is Standard
- Deviation

Properties of Normal Distribution

- For normal distribution of data, mean, median, and mode are equal, (i.e, Mean Median = Mode).
- Total area under the normal distribution curve is equal to 1.
- Normally distributed curve is symmetric at the center along the mean.
- In a normally distributed curve, there is exactly half value to the right of the central and exactly half value to the left side of the central value.
 - Normal distribution is defined using the values of the mean and standard deviation.

Normal distribution curve is a Unimodal Curve, i.e. a curve with only one peak how to implement in excel

1. Input your data set into an Excel spreadsheet

2. Find the mean of your data set

=AVERAGE(cell range)

"cell range" is a required component and the range of cells where your data exists, such as cells A1 through A64. You can write this in the function as A1:A64.

3. Find the standard deviation of your data set

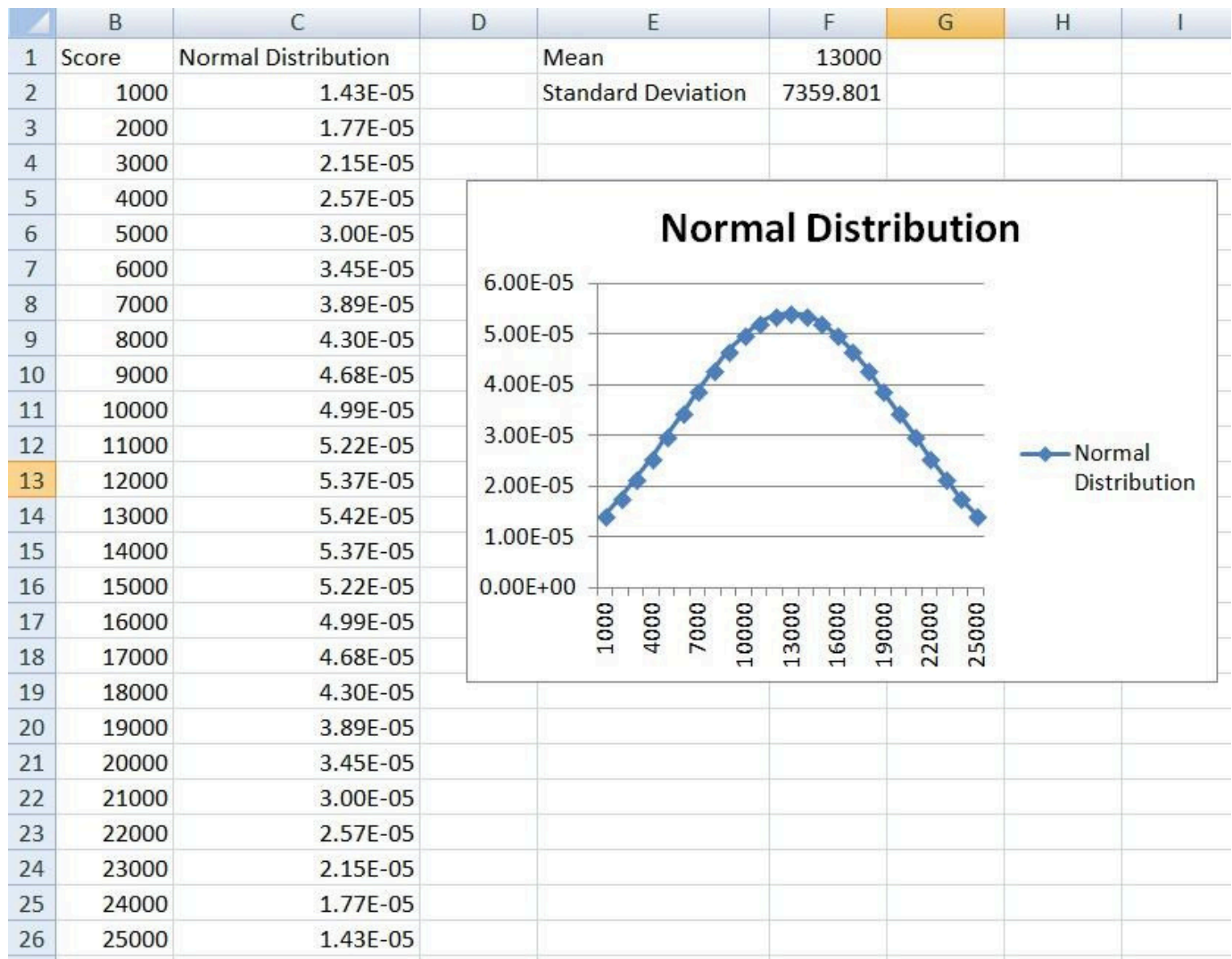
=STDEV(cell range)

4. Select a value for the distribution

5. Type the NORM.DIST function and fill

NORM.DIST(x,mean, standarddeviation,cumulative)

NORM.DIST(x,mean, standarddeviation, FALSE)



8. Calculation of cumulative distribution functions for Exponential and Normal distribution.

	A	B	C	D	E
1			NORMAL DISTRIBUTION		EXPONENTIAL DISTRIBUTION
2	Employee	Incentives	Normal Distribution	No. of days	Exponential Distribution
3	EMP1	1000	0.06859975	300	0.527633447
4	EMP2	2000	0.123838134	320	0.550671036
5	EMP3	3000	0.204480671	350	0.58313798
6	EMP4	4000	0.310147008	400	0.632120559
7	EMP5	5000	0.434415098	420	0.650062251
8	EMP6	6000	0.565584902	450	0.675347533
9	EMP7	7000	0.689852992	500	0.713495203
10	EMP8	8000	0.795519329	520	0.727468207
11	EMP9	9000	0.876161866	550	0.747160404
12	EMP10	10000	0.93140025	600	0.77686984
13					
14					
15	mean	5500		μ	400
16	standard d	3027.6504		λ	0.0025

9. Given data from two distributions, find the distance between the distributions.

Euclidean distance

Euclidean distance is the distance between two real distinct value. It is calculated by the square root of the sum of the squared difference elements in two vectors.

$$\text{Euclidean Distance} = |X - Y| = \sqrt{\sum_{i=1}^{i=n} (x_i - y_i)^2}$$

X: Array or vector X

Y: Array or vector Y

x_i: Values of horizontal axis in the coordinate plane

y_i: Values of vertical axis in the coordinate plane

n: Number of observations

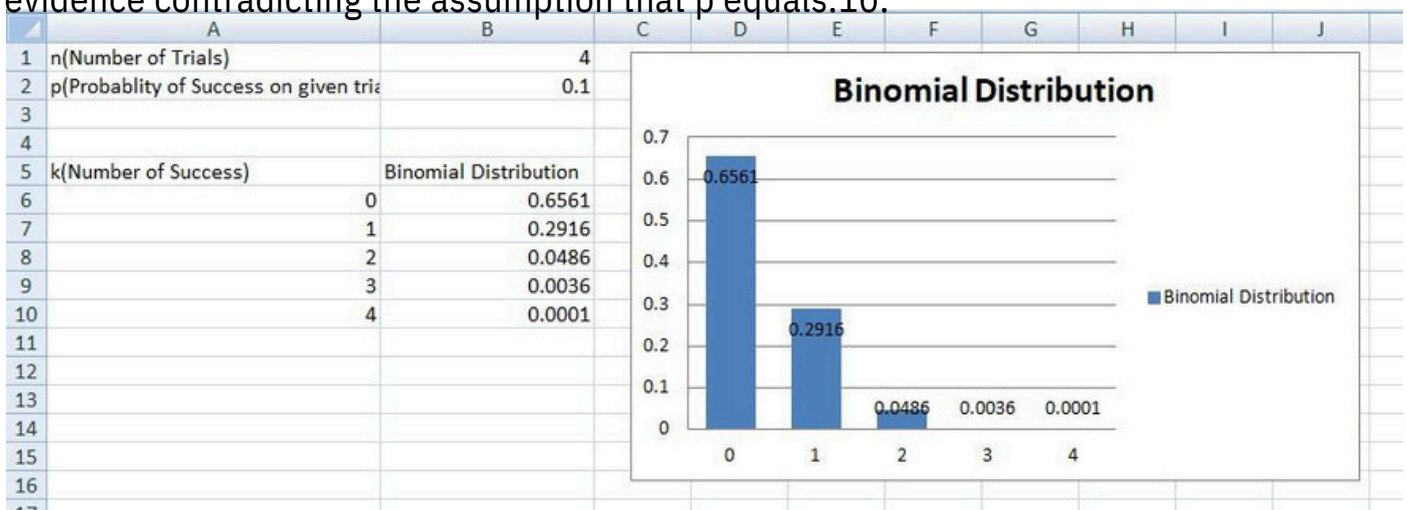
how to implement in excel

= SORT(SUM X MYZ(array_X, array_Y))

	A	B	C	D	E	F	G
1	Trials	Binomial Distribution	Poisson Distribution	Distance			
2	1	0.08957952			0.178548	n	0.4
3	2	0.20901888			0.155393	k	8
4	3	0.27869184			0.271542	λ	0.4
5	4	0.2322432			0.231528		
6	5	0.12386304			0.123806		
7	6	0.04128768			0.041284		
8	7	0.00786432			0.007864		
9	8	0.00065536			0.000655		
10							
11							

10. Application problems based on the Binomial distribution.

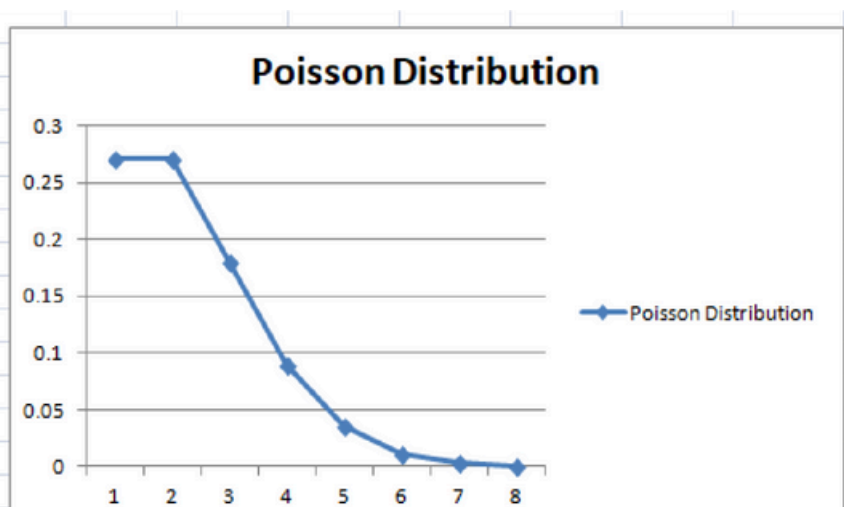
Ques: Antibiotics occasionally cause nausea as a side effect. A major drug company has developed a new antibiotic called Phe-Mycin. The company claims that, at most, 10 percent of all patients treated with Phe-Mycin would experience nausea as a side effect of taking the drug. Suppose that we randomly select $n = 4$ patients and treat them with Phe-Mycin. Each patient will either experience nausea (which we arbitrarily call a success) or will not experience nausea (a failure). We will assume that p , the true probability that a patient will experience nausea as a side effect, is .10, the maximum value of p claimed by the drug company. Furthermore, it is reasonable to assume that patients' reactions to the drug would be independent of each other. Let x denote the number of patients among the four who will experience nausea as a side effect. It follows that x is a binomial random variable, which can take on any of the potential values 0, 1, 2, 3, or 4. That is, anywhere between none of the patients and all four of the patients could potentially experience nausea as a side effect. Suppose that we wish to investigate whether p , the probability that a patient will experience nausea as a side effect of taking Phe-Mycin, is greater than .10, the maximum value of p claimed by the drug company. This assessment will be made by assuming, for the sake of argument, that p equals .10, and by using sample information to weigh the evidence against this assumption and in favor of the conclusion that p is greater than .10. Suppose that when a sample of $n=4$ randomly selected patients is treated with PheMycin, three of the four patients experience nausea. Because the fraction of patients in the sample that experience nausea is $3/4 = 0.75$, which is far greater than .10, we have some evidence contradicting the assumption that p equals .10.



11. Application problems based on the Poisson distribution.

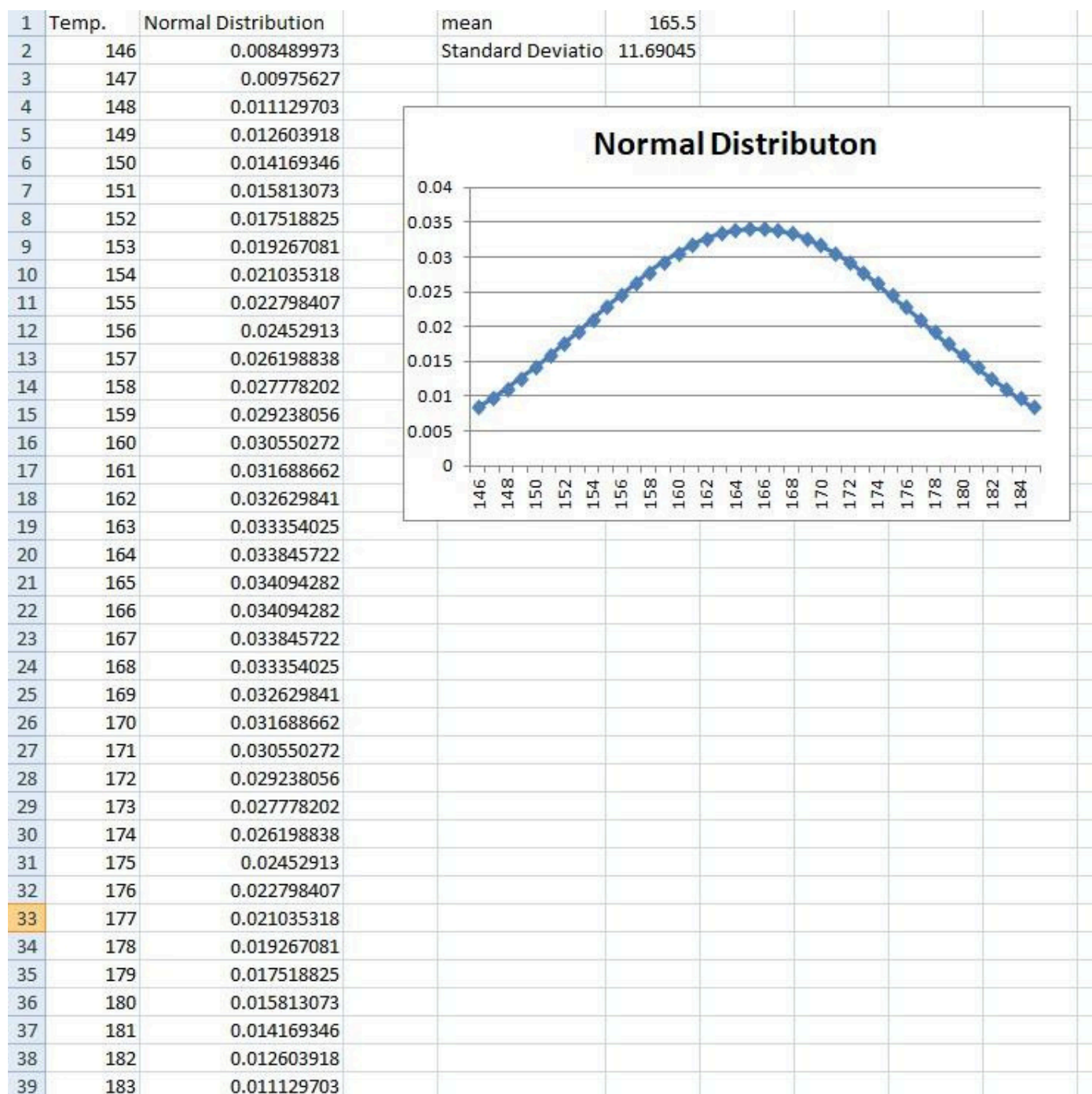
Ques: In a cafe, the customer arrives at a mean rate of 2 per min. Find the probability of arrival of 5 customers in 1 minute using the Poisson distribution formula.

1	λ	2
2	x	5
3		
4		
5	x	Poisson Distribution
6	1	0.270670566
7	2	0.270670566
8	3	0.180447044
9	4	0.090223522
10	5	0.036089409
11	6	0.012029803
12	7	0.003437087
13	8	0.000859272
14		
15		



12.Application problems based on the Normal distribution.

Ques:: According to the website of the American Association for Justice, ^11 Stella Liebeck of Albuquerque, New Mexico, was severely burned by McDonald's coffee in February 1992. Liebeck, who received third-degree burns over 6 percent of her body, was awarded \$160,000 in compensatory damages and \$480,000 in punitive damages. A post-verdict investigation revealed that the coffee temperature at the local Albuquerque McDonald's had dropped from about 185 degree F before the trial to about 158 degree after the trial. This case concerns coffee temperatures at a fast-food restaurant. Because of the possibility of future litigation and to possibly improve the coffee's taste, the restaurant wishes to study the temperature of the coffee it serves. To do this, the restaurant personnel measure the temperature of the coffee being dispensed (in degrees Fahrenheit) at a randomly selected time during each of the 24 half-hour periods from 8 a.m. to 7:30 p.m on a given day. This is then repeated on a second day, giving the 48 coffee temperatures in excel.



13. Presentation of bivariate data through scatter-plot diagrams and calculations of covariance.

Bivariate Data/ Bivariate Analysis

Bivariate analysis is one of the statistical analysis where two variables are observed. One variable here is dependent while the other is independent. These variables are usually denoted by X and Y. So, here we analyse the changes occurred between the two variables and to what extent.

The term bivariate analysis refers to as the analysis of two variables, the objective of bivariate analysis to understand the relationship between two variables. There are three common way to analysis the bivariate analysis -

1. Scatter plots
2. Correlation Coefficient
3. Simple linear Regression(SLR)

Bivariate frequency distribution

A series of statistical data showing the frequency of two variables simultaneously is called Bivariate frequency distribution. In other words, the frequency distribution of two variable is called Bivariate frequency distribution. For example: sales and advertisement expenditure, weight and height of an individual.

Why bivariate frequency distribution is significant in business research ?

1. Decision Making
2. Market-segmentation
3. Risk-assessment
4. Resource allocation how to implement in excel

= COVARIANCE.P(array1,array2)

The COVARIANCE.P function used the following arguments array 1, this is range or array of integer value.

array2 is also the second range or values.

Few things to remember about argument

1. If the given array contain text or logical value then are ignore by the Covariance function in excel.
2. The data should contain numbers, names, array or references that are numeric. IF the some cell donot contain numeric data they are ignored.
3. The data set should be same size with the same number of data points.
4. The data set should not be empty nor should the standard Deviation of the value equal.

$$\text{cov}(X, Y) = \frac{\sum (X_i - \bar{X}) * (Y_i - \bar{Y})}{n}$$

X and Y are the sample mean of the two set of values and n is the sample size.

5. Covariance is measure to indicate the extent to which two random variable in tandem.

6. Correlation is the measure used to represent how strongly two random variable are strongly related to each other.

7. Covariance is nothing but a measure of correlation.

8. Correlation referred to the scaled form of covariance.

9. Covariance can vary between -∞ to +∞ and correlation range between -1 to +1.

10. Covariance indicate the direction of the linear relationship between variables.

11. Correlation on the other hand measure both the strength and direction of the linear relationship between two variables.

12. Covariance is affected by change in scale.

13. Correlation is not affected by the change in scale.

Pearson Correlation Coefficient formula

$$r = \frac{\sum (x_i - \bar{x}) (y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

r = correlation coefficient

x_i = values of the x-variable in a sample

\bar{x} = mean of the values of the x-variable

y_i = values of the y-variable in a sample

\bar{y} = mean of the values of the y-variable



Positive correlation

As one variable increases so does the other variable.



Negative correlation

As one variable increases the other variable decreases.



No correlation

There is no relationship between the two variables.

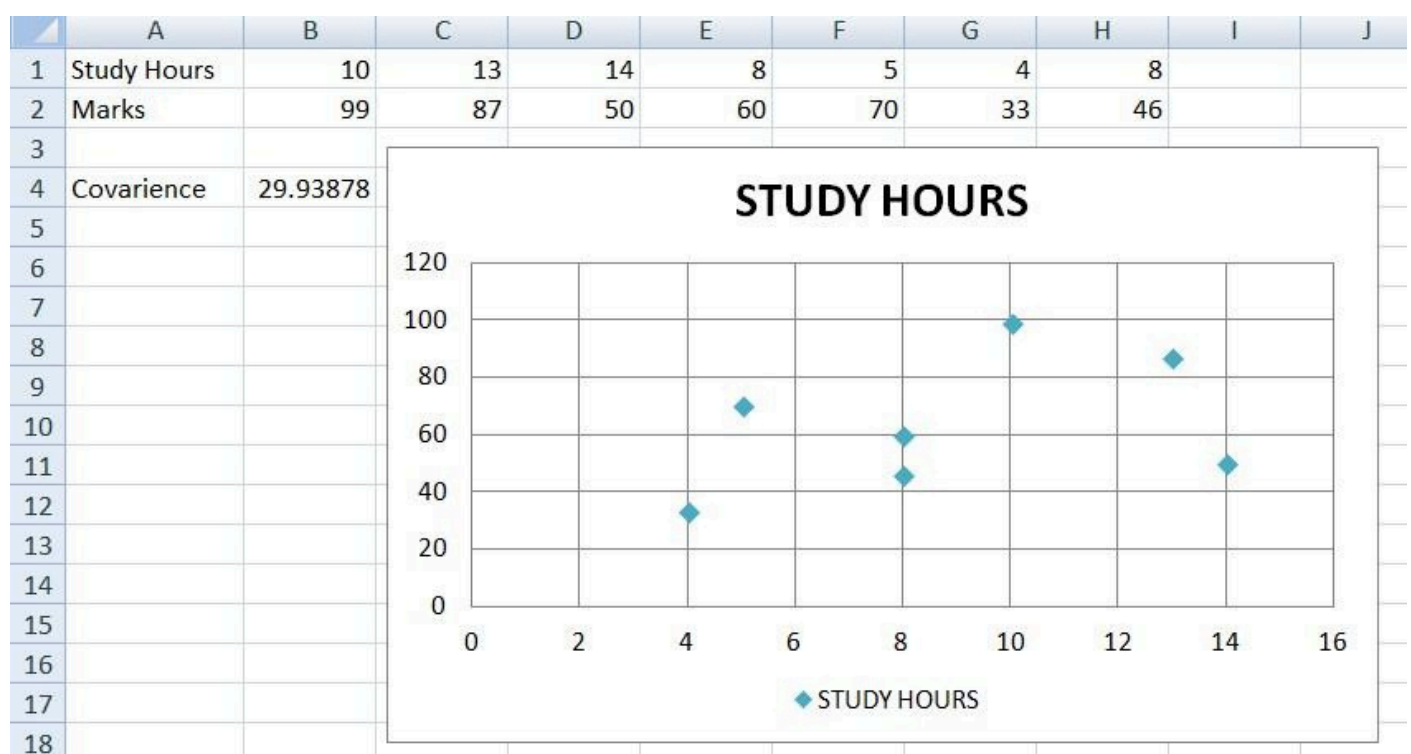
= PEARSON(array:array2)

Scatter plots:

Scatter plots are the graphs that present the relationship between two variables in a data-set. It represents data points on a two-dimensional plane or on a Cartesian system. The independent variable or attribute is plotted on the X-axis, while the dependent variable is plotted on the Y-axis. These plots are often called scatter graphs or scatter diagrams.

Scatter plots instantly report a large volume of data. It is beneficial in the following situations

- For a large set of data points given
- Each set comprises a pair of values
- The given data is in numeric form



Click insert tab along the top ribbon then click scatter chart within chart group.

CORRELATION in excel and COVARIANCE in excel -

= CORREL(hours,score)

= COVARIANCE.P(hours, score)

14. Calculation of Karl Pearson's correlation coefficients.

	A	B	C	D	E
1	X	Y1	Y2	Y3	
2	2	80	65	10	
3	5	95	69	30	
4	6	76	60	15	
5	8	58	95	25	
6	10	67	80	10	
7					
8	Karl Pearson's Correlation	XY1	-0.6273178		
9		XY2	0.63322478		
10		XY3	0.01814885		
11					

15. To find the correlation coefficient for a bivariate frequency distribution.

	A	B	C	D	E	F	G	H	I
1			Age in Years					Marginal Frequency Distribution of X:	
2	Marks	16_18	18_20	20_22	22_24	Total		Marks	Total
3	10_20	2	1	1	0	4			15 4
4	20_30	3	2	3	1	9			25 9
5	30_40	3	3	5	6	17			35 17
6	40_50	2	2	3	4	11			45 11
7	50_60	0	1	2	2	5			55 5
8	60_70	0	1	2	1	4			65 4
9									
10								Age in Years	Total
11	correlation coefficient along x-axis								17 10
12	-0.18773								19 10
13									21 16
14	correlation coefficient along Y-axis								23 14
15	0.774597								

16 .Generating Random numbers from discrete (Bernoulli, Binomial, Poisson) distributions.

How to implement in excel-

= BINOM.INV (1,P, RAND()) will generate 1 or 0 with chance of 1 being P random number

1	Random numbers from Binomial distributions				
2	N=10				
3	P=0.4				
4	5				
5	4				
6	4				
7	5				
8	4				

17 Generating Random numbers from continuous (Uniform, Normal) distributions.

= NORMINV(RAND(),B2,C2)

Where this RAND() function create your probability. B2 provides you mean, C2 refers your standard deviation.

	A	B	C
1	Random numbers from Normal distributions		
2	mean	20	
3	standard Deviation	5	
4			
5		25.03399	
6		6.594571	
7		22.82169	
8		10.33076	

18 Find the entropy from the given data set.

The entropy of a random variable is the average level of information, surprise, or uncertainty inherent to the variable's possible outcomes. Given a discrete random variable X which takes value in the alphabet \mathcal{X} and distributed according to the $P: \mathcal{X} \rightarrow [0, 1]$ The entropy is $H[X]$

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x)$$

The choice of base for log varies for different applications.

Base 2 gives the unit of bits while base e gives natural units.

Base e gives the units of $H(X)$

An equivalent definition of entropy is the expected value of the self information of a variable.

