

Overview: Use of Probability

In any statistical analysis we use probability in two ways:

1. Probability models describe non-deterministic nature of measurements.
2. Probability is used to quantify the uncertainty in the results (conclusions) of our statistical analysis.

Random Variable:

A random variable is the concept (mapping) some measurement.

- The value of the measured need not be the same every time, so we speak of the probability distribution of the random variable.

- We know everything there is to know about the random variable once we know the probability that the value of the random variable might be in any specified set.

Examples:

Sex
Age
Height
Blood Pressure
Pulse beat rate
etc

Probability Density Function (p.d.f.):

Everything that can be known about a random variable is contained in its probability density function.

Probability density function is a rule that allows us to determine the probability that a particular measurement of a random variable might be within some set of values.

Overview: Use of Probability

In any statistical analysis we use probability in two ways:

1. Probability models describe nondeterministic nature of measurements.
2. Probability is used to quantify the uncertainty in the results (conclusions) of our statistical analysis.

Random Variable:

A random variable is a function on sample space.

It is basically a device for transferring probabilities from complicated sample spaces to simple sample spaces.

A random variable X is a function whose domain is the sample space and whose range is the set of real numbers.

Thus a random variable assigns a real value (i.e., a number) to every outcome in the sample space. The particular values are called realisations and are denoted as x .

If the realisations are countable, x_1, x_2, \dots , the random variable is said to be discrete.

In contrast, if there are infinite-many uncountable realisations, the random variable is said to be continuous.

Thus, random variables can be classified according to the types of values they can take on (the range of the random variables).

Binary, categorical (ordered or unordered)
Qualitative, Quantitative.

Random Variable:

A random variable is a function on the sample space. It is basically a device for transferring probabilities from complicated sample spaces to simple sample spaces where the elements are just natural numbers.

Suppose the arrival of telephone calls at an exchange. Modelling this is very complicated as the sample space should include all possible times of arrival of calls and all possible number of calls. If we consider the random variable which counts how many calls arrive before time t (for example) then the sample space becomes $S = \{0, 1, 2, \dots\}$.

Thus a random variable X is a function whose domain is the sample space and whose range is the set of real numbers. The random variable assigns a real value (i.e. a number) to every outcome in the sample space. The particular values are called realizations and are denoted by x .

If the realizations are countable, x_1, x_2, \dots , the random variable is said to be discrete.

In contrast, if there are infinitely-many uncountable realizations, the random variable is said to be continuous.

Probability Mass Function!

If X is a discrete random variable, the function given by $f(x) = P(X=x)$ for every x within the range of X , is called Probability mass function or probability function or probability distribution of X .

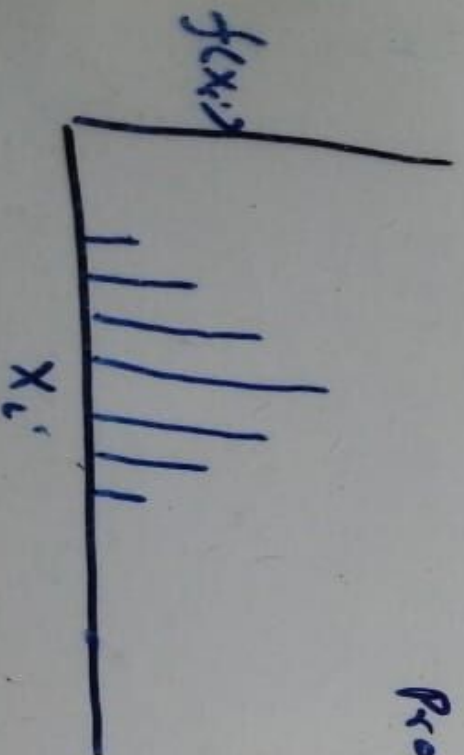
A function can serve as the probability distribution of a discrete random variable X iff its values, $f(x)$, satisfy the conditions.

1. $f(x) \geq 0$ for all x within its range of X .
2. $\sum_x f(x) = 1$, where the summation extends over the values within its domain.

Thus a table, formula or graph representing the probabilities that a random variable X takes, is called a probability distribution.

x_i	x_1	x_2	...	x_k	Total
$f(x_i)$	$f(x_1)$	$f(x_2)$...	$f(x_k)$	1

Probability Mass Function.



Probability Distribution: Continuous r.v.

A random variable 'X' is said to be continuous if there exist a non-negative function f defined for all real $x \in (-\infty, \infty)$, having the property that for any set of real numbers

$$P\{X \in A\} = \int_A f(x) dx$$

The function f is called the probability density function of the random variable X . Since X must assign some value, f must satisfy

$$1. P\{X \in (-\infty, \infty)\} = \int_{-\infty}^{\infty} f(x) dx = 1$$

If $A = \{a, b\}$, Then

$$P\{X \in (a, b)\} = \int_a^b f(x) dx = 1$$

$$2. P\{X = a\} = \int_a^a f(x) dx = 0$$

$$3. P\{a \leq X \leq b\} = \int_a^b f(x) dx = ?$$

Hence for a continuous random variable

$$P\{X \leq a\} = P\{X \leq a\} =$$

$$= F(a) = \int_{-\infty}^a f(x) dx$$



is called the cumulative density function or distribution function.

Cumulative Distribution Function (CDF)

The cumulative distribution function (CDF) $F(x)$ of a discrete random variable X with p.m.f $P(x)$ or $f(x)$ is defined for every number x by

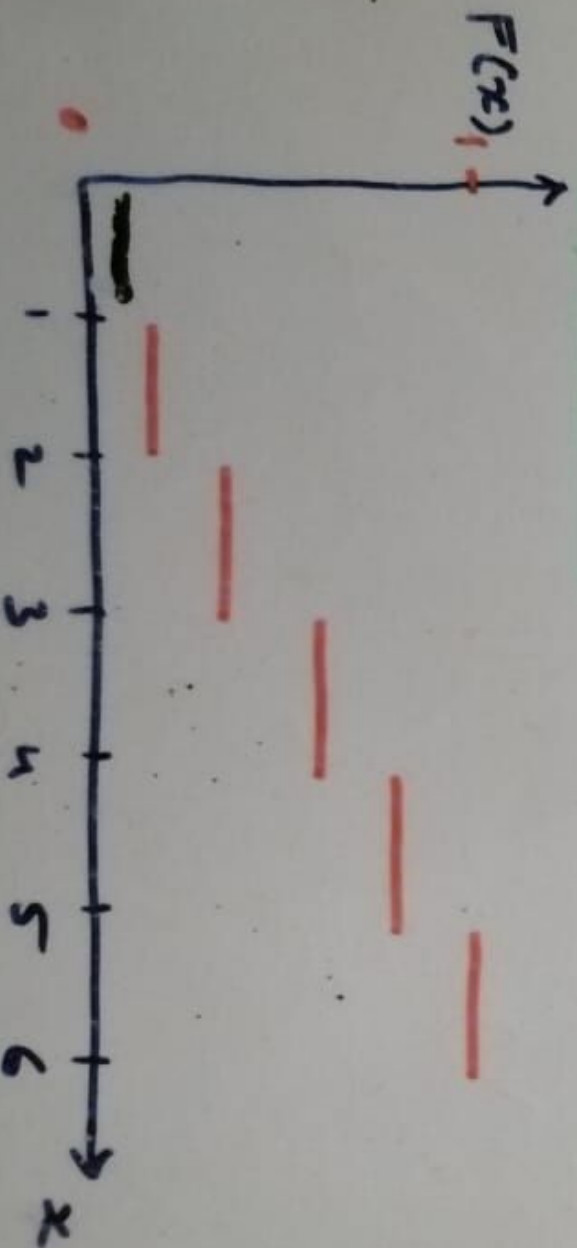
$$F(x) = P(X \leq x) = \sum_{y: y \leq x} P(y) = \sum_{x=0}^x P(x)$$

For any number x , $F(x)$ is the probability that the observed value of X will be at most x ,

For a continuous random variable

$$F(a) = P\{X \leq a\} = P\{X \leq a\} = \int_{-\infty}^a f(x) dx$$

For X a discrete r.v., the graph will have a jump at every possible value of X and will be flat between any two possible values of X . Such a graph is called a step function.



Cumulative Distribution Function

The Cumulative distribution function (cdf) $F(x)$ of a discrete random variable X with p.m.f $P(x)$ is defined for every number x by

$$a) \quad F(x) = P(X \leq x) = \sum_{y: y \leq x} P(y)$$

For any number x , $F(x)$ is the probability that an observed value of X will be at most x .

Example:

Suppose a p.m.f of Y is

$$Y \quad 1 \quad 2 \quad 3 \quad 4$$

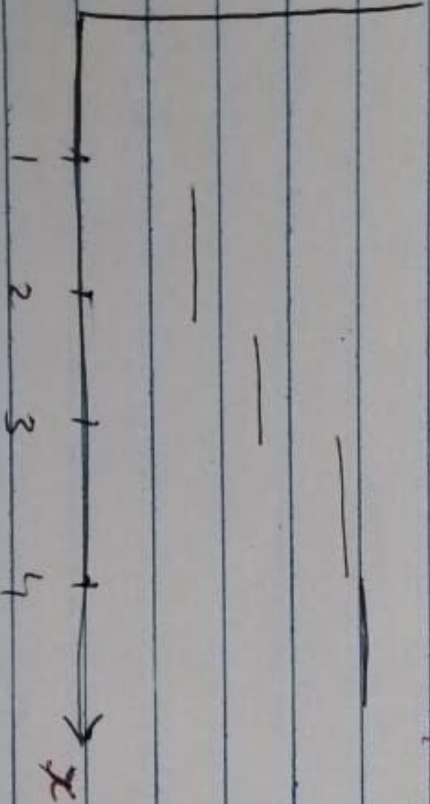
$$P(Y) \quad .4 \quad .3 \quad .2 \quad .1$$

$$F(Y) = \begin{cases} 0 & \text{if } Y < 1 \\ 0.4 & \text{if } 1 \leq Y < 2 \\ 0.7 & \text{if } 2 \leq Y < 3 \\ .9 & \text{if } 3 \leq Y < 4 \\ 1 & \text{if } 4 \leq Y \end{cases}$$

A graph of $F(x)$ is (b) Hence for a continuous r.v.

$$F(a) = P\{X \leq a\} = P\{X \leq a\} \\ = \int_{-\infty}^a f(x) dx.$$

$F(x)$



For X a d.r.v., the graph of $F(x)$ will have a jump at every possible value of X and will be flat between possible values. Such a graph is called a step-function.

Problem:

A gas station operates two pumps, each of which can pump up to 10,000 gallons of gas in a month. The total amount of gas pumped at the station in a month is a random variable Y (measured in 10,000 gallons) with a p.d.f. given by

$$f(y) = \begin{cases} y & , 0 < y < 1 \\ 2-y & , 1 \leq y < 2 \\ 0 & \text{elsewhere} \end{cases}$$

- Graph $f(y)$
- Find $F(y)$ and graph it
- Find the probability that the station will pump between 8,000 and 12,000 gallons in a particular month.
- Given that the station pumped more than 10,000 gallons in a particular month, find the probability that the station pumped more than 15,000 gallons during the month.

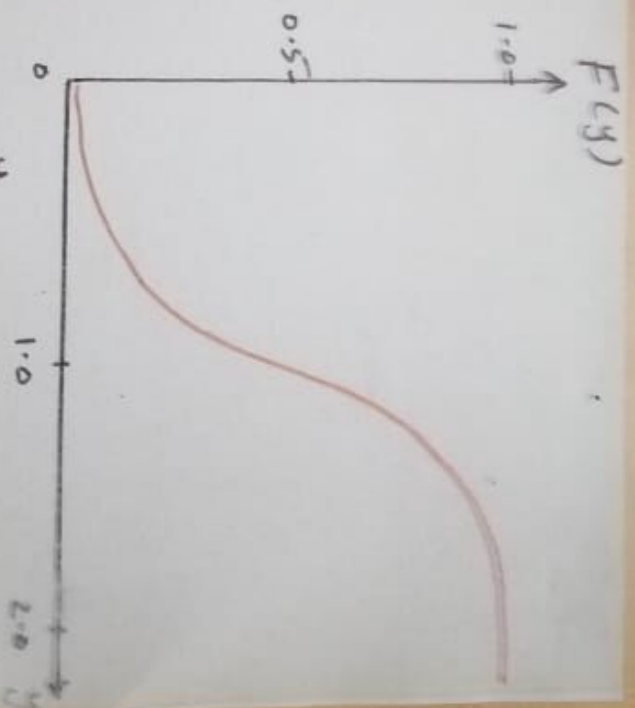
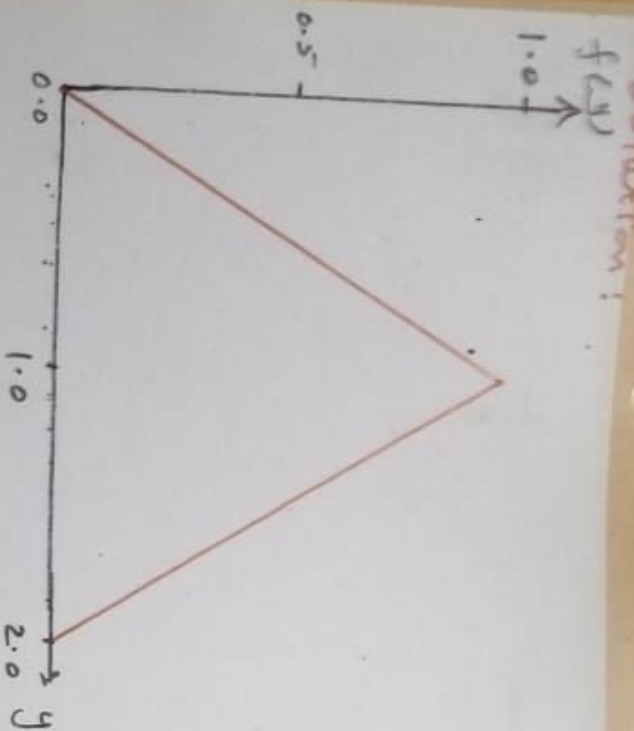
Problem:

A gas station operates two pumps, each of which can pump up to 10,000 gallons of gas in a month. The total amount of gas pumped at the station in a month is a random variable Y (measured in 10,000 gallons) with a p.d.f given by

$$f(y) = \begin{cases} y & , 0 < y < 1 \\ 2-y & , 1 \leq y < 2 \\ 0 & \text{elsewhere} \end{cases}$$

- Graph $f(y)$
- Find $F(y)$ and graph it
- Find the probability that the station will pump between 8,000 and 12,000 gallons in a particular month.
- Given that the station pumped more than 10,000 gallons in a particular month, find the probability that the station pumped more than 15,000 gallons during the month.

Solution:



$$b) F(y) = P(Y \leq y) = \int_{-\infty}^y f(y) dy$$

$$= \begin{cases} \int_0^y y dy & ; 0 < y < 1 \\ \int_0^1 y dy + \int_1^y (2-y) dy & ; 1 \leq y < 2 \\ 1 & ; y \geq 2 \end{cases}$$

$$= \begin{cases} \frac{y^2}{2} & ; 0 < y < 1 \\ \frac{y^2}{2} \Big|_0^1 + 2y \Big|_1^y - \frac{y^2}{2} \Big|_1^y & ; 1 \leq y < 2 \\ 1 & ; y \geq 2 \end{cases}$$

$$= \begin{cases} \frac{y^2}{2} & ; 0 < y < 1 \\ \frac{1}{2} + 2(y-1) - \frac{1}{2}(y^2-1) & ; 1 \leq y < 2 \\ 1 & ; y \geq 2 \end{cases}$$

$$F(y) = \begin{cases} \frac{y^2}{2} & ; 0 < y < 1 \\ 2y - \frac{y^2}{2} - 1 & ; 1 \leq y < 2 \\ 1 & ; y \geq 2 \end{cases}$$

$$\begin{aligned} c) P(8000 \leq y \leq 12000) &= P(0.8 \leq y \leq 1.2) \\ &= \int_{0.8}^1 y \, dy + \int_1^{1.2} (2-y) \, dy \\ &= \frac{1}{2} [y^2]_{0.8}^1 + 2[y]_1^{1.2} - \frac{1}{2} [y^2]_1^{1.2} \\ &= \frac{1}{2} (1 - 0.64) + 2(1.2 - 1) - \frac{1}{2} (1.44 - 1) \\ &= 0.36 \end{aligned}$$

$$\begin{aligned} d) P(Y > 15000 \cap X > 10000) &= \frac{P(Y > 15000)}{P(X > 10000)} \\ &= \frac{P(Y > 1.5)}{P(X > 1.0)} = \frac{1 - P(Y \leq 1.5)}{1 - P(X \leq 1.0)} \\ P(Y \leq 1.5) &= \int_0^1 y \, dy + \int_1^{1.5} (2-y) \, dy \\ &= \frac{1}{2} [y^2]_0^1 + 2[y]_1^{1.5} - \frac{1}{2} [y^2]_1^{1.5} \\ &= \frac{1}{2} (1 - 0) + 2(1.5 - 1) - \frac{1}{2} (2.25 - 1) \\ &= 0.875 \end{aligned}$$

PRACTICAL 2: Probability Distributions

1. A random variable 'X' has a probability distribution.

$$P(X) = \begin{cases} k/x & ; x=2,4,8 \\ 0 & ; \text{otherwise} \end{cases}$$

- a) Determine the value of the constant k
 b) what is the value of $P[X=4]$?
 c) what is $P[X < 4]$?
 d) what is $P[3 \leq X \leq 9]$?

2. A random variable 'X' has a CDF

$$F(X) = \begin{cases} 0 & X < -3 \\ 0.4 & -3 \leq X < 5 \\ 0.8 & 5 \leq X < 7 \\ 1 & X \geq 7 \end{cases}$$

- a) Draw a graph of the CDF
 b) write $P(X)$, the P.d.f of X .

3. Sketch the following functions are p.d.f for some values of k and determine k . Also sketch the functions

$$\begin{aligned} \text{a) } f(X) &= kx^2 & ; & 0 \leq X \leq 2 \\ \text{b) } f(X) &= k(2-X) & ; & 0 \leq X \leq 1 \\ \text{c) } f(X) &= ke^{-x} & ; & 0 < X \end{aligned}$$

4. A continuous random variable has p.d.f

$$f(X) = kx \quad ; \quad 0 \leq X \leq 4$$

- a) Find the value of the constant k .
 b) sketch the function $f(X)$
 c) find $P[1 \leq X \leq 2.5]$ and sketch

Assignment 2

Q1. (6 points) The sample space of a random experiment is $\{a, b, c, d, e, f\}$ and each outcome is equally likely. A random variable is defined as follows

outcome	x		x		x	
	a	b	c	d	e	f
	0	0	1.5	1.5	2	3

Determine the probability mass function of X . Determine the following probabilities:

- (a) $P(X = 1.5)$ (b) $P(0.5 < X < 2.7)$ (c) $P(X > 3)$
 (d) $P(0 \leq X < 2)$ (e) $P(X = 0 \text{ or } X = 2)$

Solution to Q1:

Probability mass function is

$$P(X = 0) = P(\{a, b\}) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6}; \quad P(X = 1.5) = \frac{2}{6}, \quad P(X = 2) = \frac{1}{6}, \quad P(X = 3) = \frac{1}{6}.$$

- (a) $P(X = 1.5) = \frac{2}{6}$ (b) $P(0.5 < X < 2.7) = P(X = 1.5) + P(X = 2) = \frac{3}{6}$
 (c) $P(X > 3) = 0$ (d) $P(0 \leq X < 2) = P(X = 0) + P(X = 1.5) = \frac{4}{6}$
 (e) $P(X = 0 \text{ or } X = 2) = \frac{3}{6}$

Marking scheme for Q1:

Completely correct p.m.f. - 1 point, correct answer for each part - 1 point. Total - 6 points.

Q2. (2 points) Determine the mean and the variance in Question Q1

Solution to Q2:

$$E(X) = 1.33, \quad \text{Var}(X) \approx 1.15$$

Marking scheme for Q2:

1 point for the mean and the variance. Total - 2 points.

Q3. We say that X has *uniform distribution* on a set of values $\{x_1, \dots, x_k\}$ if

$$P(X = x_i) = \frac{1}{k}, \quad i = 1, \dots, k.$$

The thickness measurements of a coating process are *uniformly distributed* with values 0.15, 0.16, 0.17, 0.18, 0.19. Determine the mean and variance.

Solution to Q3:

We have $P(X = 0.15) = \dots = P(X = 0.19) = 1/5$. Mean: 0.17; Variance:

$$\frac{1}{5} (0.01^2 + 0.02^2 + 0.02^2 + 0.01^2)$$

Marking scheme for Q3:

This question will not be marked

2

Binomial Experiment :

A binomial experiment is one that possess the following properties.

- a) The experiment consists of n repeated trials
- b) Each trial results in an outcomes that may be classified as success or a failure.
- c) The probability of success, denoted by p remains constant from trial to trial.
- d) The repeated trials are independent.

Binomial Distribution:

Consider an experiment with two possible outcomes call them Success (S) and failure (f) with $P(S) = p$ and $P(f) = q$ such that $p+q=1$. Let 'X' be a random variable denotes the number of successes in n independent repeated trials, e.g.

consider $\overbrace{S \ S \ \dots \ S \ \ f \ f \ \dots \ f}^{n \text{ trials}}$
 $\underbrace{\hspace{1.5cm}}_{x \text{ Success}} \quad \underbrace{\hspace{1.5cm}}_{n-x \text{ failure}}$

The probability distribution of the particular sequence (by multiplicative Law of independent events) is

$$p^x q^{n-x} \quad \text{or} \quad p^x (1-p)^{n-x}$$

The number of sequence in which 'x' successes and $n-x$ failures are observed in some order is $\binom{n}{x}$ ways. which is binomial coefficient.

Thus the probability Distribution that exactly x successes and $n-x$ failures occur in n independent trials is

$$b(x; n, p) = \binom{n}{x} p^x q^{n-x}; \quad x=0, 1, 2, \dots, n$$

which is known as binomial distribution with index 'n' and parameter p.

$$\mu'_4 = E[X^4] = \sum_{x=0}^n x^4 \binom{n}{x} p^x q^{n-x}$$

$$= \sum_{x=0}^n [x(x-1)(x-2)(x-3) + 6x(x-1)(x-2) + 7x(x-1) + x] \binom{n}{x} p^x q^{n-x}$$

$$\mu'_4 = n(n-1)(n-2)(n-3)p^4 + 6n(n-1)(n-2)p^3 + 7n(n-1)p^2 + np$$

Moment about the mean:

$$\mu_r = \frac{1}{n} \sum (x_i - \bar{x})^r$$

$$\mu_1 = \mu'_1 - \mu'_1 = 0$$

$$\mu_2 = \mu'_2 - (\mu'_1)^2$$

$$= n(n-1)p^2 + np - n^2p^2$$

$$= n^2p^2 - np^2 + np - n^2p^2$$

$$= np(1-p)$$

$$\sigma^2 = \boxed{\mu_2 = npq} \Rightarrow \sigma = \sqrt{\mu_2} = \sqrt{npq}$$

$$\mu_3 = \mu'_3 - 3\mu'_2\mu'_1 + 2\mu'_1{}^3$$

$$= [n(n-1)(n-2)p^3 + 3n(n-1)p^2 + np]$$

$$- 3[n(n-1)p^2 + np]np + 2n^3p^3$$

$$= n[n^2 - 3n + 2]p^3 + 3[n^2 - n]p^2 + np - 3n^2p^2 - 3n^2p^3(n-1) + 2n^3p^3$$

$$= n^3p^3 - 3n^2p^3 + 2np^3 + 3n^2p^2 - 3np^2 + np - 3n^2p^2 - 3n^3p^3 + 3n^2p^3 + 2n^3p^3$$

$$= 2np^3 - 3np^2 + np = np[1 - 3p + 2p^2]$$

$$\mu'_1 = nP \left[q^{\binom{n}{1}} + \binom{n-1}{1} P q^{n-2} + \dots + P^{n-1} \right]$$

$$= nP [q + P]^{n-1}$$

$$\boxed{\mu'_1 = nP}$$

$$\therefore q + P = 1$$

at $r=2$

$$\mu'_2 = E[X^2] = \sum_{X=0}^n X^2 \binom{n}{X} P^X q^{n-X}$$

$$= \sum_{X=0}^n [X + X(X-1)] \binom{n}{X} P^X q^{n-X}$$

$$= \sum_{X=0}^n X \binom{n}{X} P^X q^{n-X} + \sum_{X=0}^n X(X-1) \binom{n}{X} P^X q^{n-X}$$

$$= nP + 0 + 0 + 2 \cdot 1 \binom{n}{2} P^2 q^{n-2} + 3 \cdot 2 \binom{n}{3} P^3 q^{n-3} + \dots + n(n-1) \binom{n}{n} P^n q^{n-n}$$

$$= nP + n(n-1) P^2 [q^{n-2} + \binom{n-2}{1} P q^{n-3} + \dots + P^{n-2}]$$

$$= nP + n(n-1) P^2 [q + P]^{n-2}$$

$$\checkmark = nP + n(n-1) P^2$$

$$= nP [1 + (n-1)P] = nP [1 + nP - P]$$

$$= nP [q + nP]$$

$$\boxed{\mu'_2 = nPq + n^2 P^2}$$

Similarly Proceeding we get

$$\mu'_3 = E[X^3] = \sum_{X=0}^n X^3 \binom{n}{X} P^X q^{n-X}$$

$$= \sum_{X=0}^n [X(X-1)(X-2) + 3X(X-1) + X] \binom{n}{X} P^X q^{n-X}$$

$$= n(n-1)(n-2) P^3 (q + P)^{n-3} + 3n(n-1) P^2 (q + P)^{n-2} + nP (q + P)^{n-1}$$

$$\boxed{\mu'_3 = n(n-1)(n-2) P^3 + 3n(n-1) P^2 + nP}$$