# Probability and Statistics ( BTech CSE )

Anmol Nawani

April 28, 2023

## Contents

## 1 Ungrouped Data

Ungrouped data is data that has not been arranged in any way.So it is just a list of observations

$$x_1, x_2, x_3, ...x_n$$

### 1.1 Mean

$$\bar{x} = \frac{x_1 + x_2 + x_3 + ... + x_n}{n}$$

$$\bar{x} = \frac{\sum_{i=1}^{n} x_i}{n}$$

### 1.2 Mode

The observation which occurs the highest number of time. So the $x_i$ which has the highest count in the observation list.

### 1.3 Median

The median is the middle most observations. After ordering the n observations in observation list in either Ascending or Descending order (any order works). The median will be :

- n is even

$$Median = \frac{x_{\frac{n}{2}} + x_{(\frac{n}{2}+1)}}{2}$$

- n is odd

$$Median = x_{\frac{n+1}{2}}$$

## 1.4 Variance and Standard Deviation

$$Variance = \sigma^2$$

$$Standard\ deviation = \sigma$$

$$\sigma^2 = \frac{\sum_{i=1}^{n}(x_i - Mean)^2}{n}$$

$$\sigma^2 = \frac{\sum_{i=1}^{n} x_i^2}{n} - (Mean)^2$$

## 1.5 Moments

### 1.5.1 About some constant A

$$r^{th}\ moment = \frac{1}{n}\Sigma(x_i - A)^r$$

### 1.5.2 About Mean (Central Moment)

When A = Mean, then the moment is called central moment.

$$\mu_r = \frac{1}{n}\Sigma(x_i - Mean)^r$$

### 1.5.3 About Zero (Raw Moment)

When A = 0, then the moment is called raw moment.

$$\mu_r^{'} = \frac{1}{n}\Sigma x_i^r$$

# 2 Grouped Data

Data which is grouped based on the frequency at which it occurs. So if 9 appears 5 times in our observations, we group as x(observation) = 9 and f (frequency) = 5.

| x (observations) | f (frequency) |
|:---:|:---:|
| 2 | 5 |
| 1 | 3 |
| 4 | 5 |
| 8 | 9 |

If we store it in data way, i.e. the observations are of form 10-20, 20-30, 30-40 … then we will get $x_i$ by doing

$$x_i = \frac{lower\ limit + upper\ limit}{2}$$

i.e,
$x_i$ for 20-30 will be $\frac{20+30}{2}$
So for data

| | f (frequency) |
|:---:|:---:|
| 0- 20 | 2 |
| 20-40 | 6 |
| 40-60 | 1 |
| 60-80 | 3 |

the $x_i$'s will become.

| | $f_i$ | $x_i$ |
|:---:|:---:|:---:|
| 0- 20 | 2 | 10 |
| 20-40 | 6 | 30 |
| 40-60 | 1 | 50 |
| 60-80 | 3 | 70 |

## 2.1 Mean

$$\bar{x} = \frac{\Sigma f_i x_i}{\Sigma f_i}$$

## 2.2 Mode

The **modal class** is the record with the row with the highest $f_i$

$$Mode = l + (\frac{f_1 - f_0}{2f_1 - f_0 - f_2}) \times h$$

In the formula :

$l \rightarrow$ lower limit of modal class
$f_1 \rightarrow$ frequency($f_i$) of the modal class
$f_0 \rightarrow$ frequency of the row preceding modal class
$f_2 \rightarrow$ frequency of the row after the modal class
$h \rightarrow$ size of class interval (upper limit - lower limit)

## 2.3 Median

The median for grouped data is calculated with the help of **cumulative frequency**. The cumulative frequency ($cf_i$) is given by:

$$cf_i = f_1 + f_2 + f_3 + ... + f_i$$

The **median class** is the class whose $cf_i$ is just greater than or is equal to $\frac{\Sigma f}{2}$

$$Median = l + (\frac{(n/2) - cf}{f}) \times h$$

In the formula :

$l \rightarrow$ lower limit of the median class
$h \rightarrow$ size of class interval (upper limit - lower limit)
$n \rightarrow$ number of observations
$cf \rightarrow$ cumulative frequency of the median class
$f \rightarrow$ frequency of the median class

## 2.4 Variance and Standard Deviation

$$Variance = \sigma^2$$

$$Standard\ deviation = \sigma$$

$$\sigma^2 = \frac{\sum_{i=1}^{n} f_i(x_i - Mean)^2}{\Sigma f_i}$$

$$\sigma^2 = \frac{\sum_{i=1}^{n} f_i x_i^2}{\Sigma f_i} - (Mean)^2$$

## 2.5 Moments

### 2.5.1 About some constant A

$$r^{th} \ moment = \frac{1}{\Sigma f_i}[\Sigma f_i(x_i - A)^r]$$

### 2.5.2 About Mean (Central Moment)

When A = Mean, then the moment is called central moment.

$$\mu_r = \frac{1}{\Sigma f_i}[\Sigma f_i(x_i - Mean)^r]$$

### 2.5.3 About Zero (Raw Moment)

When A = 0, then the moment is called raw moment.

$$\mu_r' = \frac{1}{\Sigma f_i}[\Sigma f_i x_i^r]$$

# 3 Relation between Mean, Median and Mode

$$3Median = 2Mean + Mode$$

# 4 Relation between raw and central moments

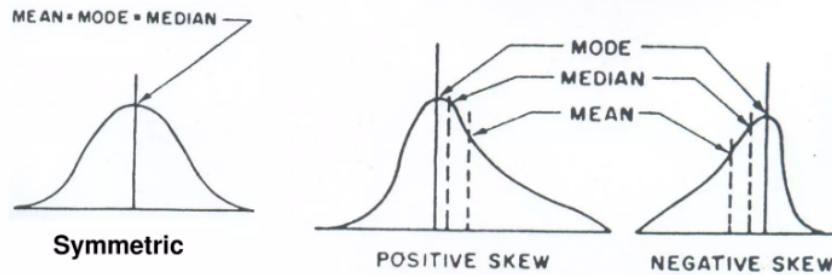$$\mu_0 = \mu_0' = 1$$
$$\mu_1 = 0$$
$$\mu_2 = \mu_2' - \mu_1'^2$$
$$\mu_3 = \mu_3' - 3\mu_1'\mu_2' + 2\mu_1'^3$$
$$\mu_4 = \mu_4' - 4\mu_3'\mu_1' + 6\mu_2'\mu_1'^2 - 3\mu_1'^4$$

# 5 Skewness and Kurtosis

## 5.1 Skewness

- If Mean > Mode, then skewness is positive
- If Mean = Mode, then skewness is zero (graph is symmetric)
- If Mean < Mode, then skewness is zero

### 5.1.1 Pearson's coefficient of skewness

The pearson's coefficient of skewness is denoted by $S_{KP}$

$$S_{KP} = \frac{Mean - Mode}{Standard\ Deviation}$$

- If $S_{KP}$ is zero then distribution is symmetrical

- If $S_{KP}$ is positive then distribution is positively skewed

- If $S_{KP}$ is negative then distribution is negatively skewed

### 5.1.2 Moment based coefficient of skewness

The moment based coefficient of skewness is denoted by $\beta_1$. The $\mu$ here is central moment.

$$\beta_1 = \frac{\mu_3^2}{\mu_2^3}$$

The drawback of using $\beta_1$ as a coefficient of skewness is that it **can only tell if distribution is symmetrical or not** ,when $\beta_1 = 0$. It can't tell us the direction of skewness, i.e positive or negative.

- If $\beta_1$ is zero, then distribution is symmetrical

### 5.1.3 Karl Pearson's $\gamma_1$

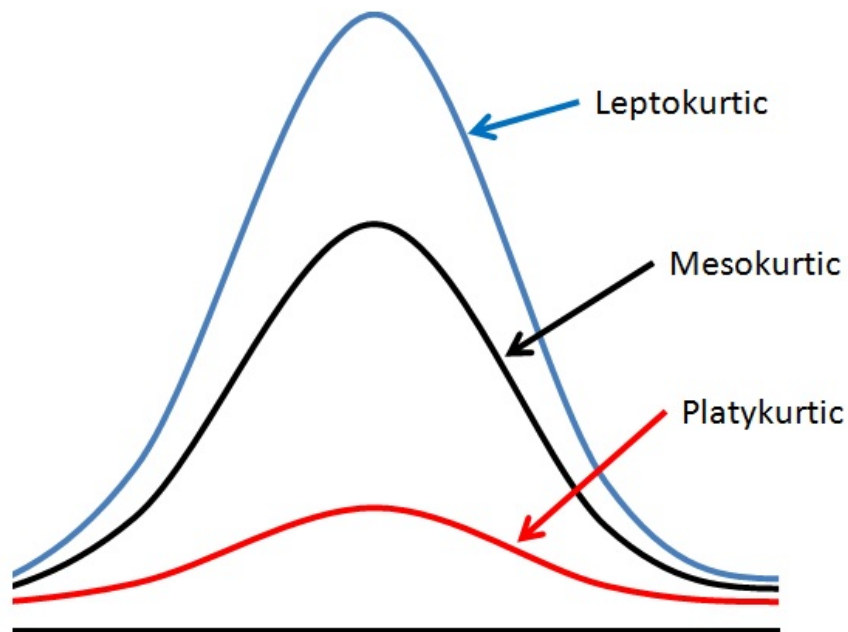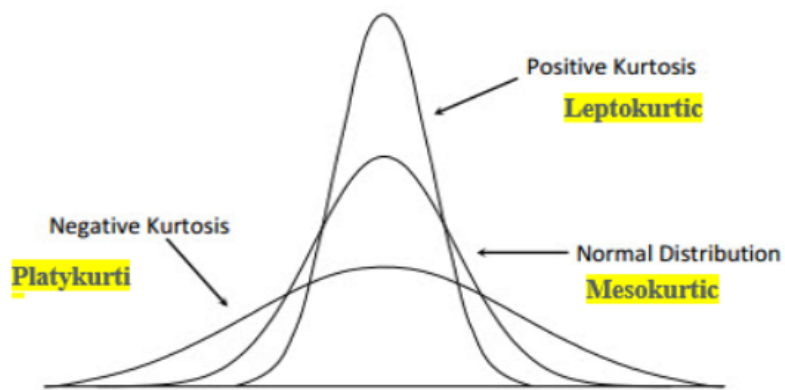To remove the drawback of the $\beta_1$ , we can derive Karl Pearson's $\gamma_1$

$$\gamma_1 = \sqrt{\beta_1}$$
$$\gamma_1 = \frac{\mu_3}{\mu_2^{3/2}}$$

6

- If $\mu_3$ is positive, the distribution has positive skewness

- If $\mu_3$ is negative, the distribution has negative skewness

- If $\mu_3$ is zero, the distribution is symmetrical

## 5.2   Kurtosis

Kurtosis is the measure of the peak and the curve and the "fatness" of the curve.

The kurtosis is calculated using $\beta_2$

$$\beta_2 = \frac{\mu_4}{\mu_2^2}$$

The value of $\beta_2$ tell's us about the type of curve

- Leptokurtic (High Peak) when $\beta_2 > 3$
- Mesokurtic (Normal Peak) when $\beta_2 = 3$
- Platykurtic (Low Peak) when $\beta_2 < 3$

### 5.2.1   Karl Pearson's $\gamma_2$

$\gamma_2$ is defined as:

$$\gamma_2 = \beta_2 - 3$$

- Leptokurtic when $\gamma_2 > 0$
- Mesokurtic when $\gamma_2 = 0$
- Platykurtic when $\gamma_2 < 0$

# 6   Basic Probability

## 6.1   Conditional Probability

If some event B has already occured, then the probability of the event A is:

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

$P(A \mid B)$ is read as A given B. So we are given that B has occured and this is probability of now A occuring.

## 6.2   Law of Total Probability

The law of total probability is used to find probability of some event A that has been partitioned into several different places/parts.

$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + P(A|B_3)P(B_3) + ... + P(A|B_i)P(B_i)$$

$$P(A) = \Sigma P(A|B_i)P(B_i)$$

**Example**, Suppose we have 2 bags with marbles

- Bag 1 : 7 red marbles and 3 green marbles

- Bag 2 : 2 red marbles and 8 green marbles

Now we select one bag at random (i.e, the probability of choosing any of the two bags is equal so 0.5). If we draw a marble, what is the probability that it is a green marble?

**Sol.** The green marbles are in parts in bag 1 and bag 2.
Let G be the event of green marble.
Let $B_1$ be the event of choosing the bag 1
Let B-2 be the event of choosing the bag 2

Then, $P(G|B_1) = \frac{3}{7+3}$ and $P(G|B_2) = \frac{8}{2+8}$
Now, we can use the law of total probability to get

$$P(G) = P(G|B_1)P(B_1) + P(G|B_2)P(B_2)$$

**Example** 2, Suppose a there are 3 forests in a park.

- Forest A occupies 50% of land and 20% plants in it are poisonous

- Forest B occupies 30% of land and 40% plants in it are poisonous

- Forest C occupies 20% of land and 70% plants in it are poisonous

What is the probability of a random plant from the park being poisonous.

**Sol.** Since probability is equal across whole area of the park. Event A is plant being from Forest A, Event B is plant being from Forest B and Event C is plant being from Forest C. If event P is plant being poisonous, then using law of total probability,

$$P(P) = P(P|A)P(A) + P(P|B)P(B) + P(P|C)P(C)$$

And we know P(A) = 0.5, P(B) = 0.3 and P(C) = 0.2. Also P(P|A) = 0.20, P(P|B) = 0.40 and P(P|C) = 0.70

## 6.3  Some basic identities

- Probabilities follow law of inclusion and exclusion

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

- DeMorgan's Theorem

$$P(\overline{A \cap B}) = P(\overline{A} \cup \overline{B})$$
$$P(\overline{A \cup B}) = P(\overline{A} \cap \overline{B})$$

- Some other Identity

$$P(\overline{A} \cap B) + P(A \cap B) = P(B)$$
$$P(A \cap \overline{B}) + P(A \cap B) = P(A)$$