**Name:** Hamza Zafar

**Email:** hamzazafar2k12@gmail.com

**Fellowship:** Bytewise- Data Science Group I

**Submitted To:** Mahrukh Khan

**Title:** Task 3

**Summary of the Assigned Topic**

# Introduction to Probability and Statistics for Data Science

**Probability and Random Variables:**

Probability measures the likelihood of an event, ranging from 0 (impossible) to 1 (certain). For instance, the probability of rolling an even number on a six-sided die is 0.5. Random variables represent outcomes of random processes; for example, the result of a dice roll is a discrete random variable taking values from 1 to 6. In contrast, continuous random variables, like bus arrival times, can take any value within a range.

**Probability Distribution:**

For discrete random variables, a probability distribution function $P(X)$ assigns probabilities to each possible value. The uniform distribution is a common discrete distribution where each outcome in the sample space has equal probability. Continuous variables require a probability density function $p(x)$, where probabilities are calculated over intervals using integrals.

**Mean, Variance, and Standard Deviation:**

The mean (average) of a sample indicates its central value. For a larger sample, this converges to the expected value $E(X)$. Variance measures the spread of data points around the mean, with the standard deviation being its square root.

**Mode, Median, and Quartiles:**

The median is the middle value separating higher and lower halves of data, while quartiles divide data into four equal parts. The mode is the most frequently occurring value. Box plots visually represent these statistics, showing the median, quartiles, and potential outliers.

**Real-World Data:**

Real-world data, like baseball players' weights, can be analysed using these statistical concepts. Data often approximates certain distributions, such as the normal distribution, where most values cluster around the mean.

**Normal Distribution:**

The normal distribution is crucial in statistics, characterized by its bell-shaped curve. It is used to model many real-world phenomena. For instance, generating random weights of potential baseball players using known mean and standard deviation produces a normal distribution.

**Confidence Intervals:**

Confidence intervals estimate population parameters (mean and variance) from sample data. They provide a range within which the true parameter likely lies with a certain confidence level. For example, the mean weight of baseball players can be estimated with different confidence levels, showing wider intervals for higher confidence.

**Hypothesis Testing:**

Hypothesis testing determines if observed differences between groups are statistically significant. For example, comparing the heights of first and second basemen involves calculating confidence intervals for their means and using the student t-test to assess the significance of the difference.

## Law of Large Numbers and Central Limit Theorem

The central limit theorem states that the mean of a large sample from any distribution will approximate a normal distribution. This theorem supports the use of normal distribution models in statistics, ensuring that with enough data, sample means converge to the population mean.

**Covariance and Correlation:**

Covariance measures the degree to which two variables change together, while correlation standardizes this measure, ranging from -1 (perfect negative correlation) to 1 (perfect positive correlation). For instance, weights and heights of baseball players show a positive correlation, indicating a relationship between the two variables.

This topic provides statistical concepts that are vital for data science: probability, distributions, statistical measures (mean, variance, etc.), hypothesis testing, and correlation. These fundamentals provide a robust foundation for analysing data, forming the basis for advanced topics and applications in data science.