

**Q1. What is the meaning of six sigma in statistics? Give proper example**

Six Sigma is a methodology used in quality management, particularly in manufacturing, to minimize defects or errors in processes. It aims to achieve near-perfect results by reducing variation and maintaining consistency. The term "Six Sigma" refers to a statistical measure of process variation where a process operates with only 3.4 defects per million opportunities (DPMO).

In statistical terms, sigma ( $\sigma$ ) represents the standard deviation, a measure of the dispersion or spread of a set of data points around the mean. The higher the sigma level, the lower the variation in the process. Six Sigma aims to bring processes to a level where they operate within six standard deviations of the mean, resulting in very high quality and minimal defects.

An example of Six Sigma implementation could be in a manufacturing process that produces automobile parts. Let's say the process involves drilling holes in metal components. Each hole must be drilled with precise dimensions, and any deviation from the specified measurements is considered a defect. By implementing Six Sigma principles, the manufacturer aims to reduce the occurrence of defective holes to no more than 3.4 per million opportunities.

To achieve this, the process undergoes rigorous analysis and improvement efforts. Statistical tools such as control charts, process capability analysis, and root cause analysis are utilized to identify and eliminate sources of variation. Quality control measures are put in place to ensure that the process consistently produces parts within the desired specifications. Through continuous monitoring and optimization, the process is brought to a Six Sigma level, resulting in exceptionally high-quality products with minimal defects.

## **Q2.What type of data does not have a log-normal distribution or a Gaussian distribution? Give proper example**

Data that do not follow a log-normal distribution or a Gaussian distribution (normal distribution) can exhibit various other types of distributions, known as non-normal distributions. Examples of such distributions include:

1. **Uniform Distribution:** In a uniform distribution, all outcomes are equally likely over a defined range. For example, when rolling a fair six-sided die, each number (1 through 6) has an equal probability of occurring.
2. **Exponential Distribution:** This distribution describes the time between events in a Poisson process, where events occur continuously and independently at a constant average rate. For instance, the time between arrivals at a bus stop follows an exponential distribution.
3. **Poisson Distribution:** It represents the number of events occurring within a fixed interval of time or space, given a constant average rate of occurrence. For instance, the number of customers arriving at a store in an hour can follow a Poisson distribution.
4. **Binomial Distribution:** This distribution describes the number of successes in a fixed number of independent Bernoulli trials, where each trial has the same probability of success. An example could be the number of heads obtained when flipping a coin multiple times.
5. **Gamma Distribution:** This distribution generalizes the exponential distribution and represents the time until the  $n$ th event in a Poisson process. It is often used in queuing theory and reliability engineering.
6. **Weibull Distribution:** It is commonly used to model the time until failure of mechanical components and represents a wide range of failure distributions.
7. **Beta Distribution:** It represents the distribution of probabilities in Bayesian statistics and is often used as a prior distribution in Bayesian inference.

These are just a few examples of distributions that do not follow a log-normal or Gaussian distribution. There are many other types of distributions, each suited to describe different types of data and phenomena.

### **Q3.What is the meaning of the five-number summary in Statistics? Give proper example**

The five-number summary is a descriptive statistic used to summarize the distribution of a dataset. It consists of five key values that divide the data into four intervals, providing insights into the center, spread, and skewness of the data. These five values are:

1. Minimum: The smallest value in the dataset.
2. First Quartile (Q1): The value below which 25% of the data fall. It represents the lower quartile.
3. Median (Second Quartile, Q2): The middle value of the dataset when it is ordered from smallest to largest. It represents the 50th percentile, or the value below which 50% of the data fall.
4. Third Quartile (Q3): The value below which 75% of the data fall. It represents the upper quartile.
5. Maximum: The largest value in the dataset.

The five-number summary is often visualized using a box plot (box-and-whisker plot), where a box is drawn from the first quartile to the third quartile, with a line representing the median. Whiskers extend from the box to the minimum and maximum values, allowing for a visual representation of the data's spread and central tendency.

Here's an example of a five-number summary for a dataset of exam scores:

Dataset: 55, 60, 65, 70, 75, 80, 85, 90, 95

1. Minimum: 55
2. First Quartile (Q1): 65
3. Median (Q2): 75
4. Third Quartile (Q3): 85
5. Maximum: 95

In this example, the first quartile (Q1) is 65, meaning 25% of the scores are below 65. The median (Q2) is 75, indicating that 50% of the scores are below 75. The

third quartile (Q3) is 85, meaning 75% of the scores are below 85. The range of the dataset is from 55 to 95, with a spread of 40 points.

#### **Q4.What is correlation? Give an example with a dataset & graphical representation on jupyter Notebook**

Correlation is a statistical measure that describes the strength and direction of a relationship between two variables. It quantifies how changes in one variable correspond to changes in another variable. Correlation values range from -1 to 1, where:

- 1 indicates a perfect positive correlation: As one variable increases, the other variable also increases linearly.
- -1 indicates a perfect negative correlation: As one variable increases, the other variable decreases linearly.
- 0 indicates no correlation: There is no apparent linear relationship between the variables.

So we create a dataset and visualize the correlation between two variables using Python in a Jupyter Notebook. First, ensure you have the necessary libraries installed (pandas, matplotlib, and seaborn). Then, follow these steps:

1. Import the required libraries.
2. Create a sample dataset with two variables.
3. Calculate the correlation coefficient.
4. Visualize the correlation using a scatter plot.