

Q-1: What are the Probability Mass Function (PMF) and Probability Density Function (PDF)? Explain with an example.

The Probability Mass Function (PMF) and Probability Density Function (PDF) are concepts used in probability and statistics to describe the distribution of random variables.

## Probability Mass Function (PMF):

The PMF is used for discrete random variables. It gives the probability that a discrete random variable is exactly equal to a specific value.

Properties:  $P(X = x) \geq 0$  for all  $x$

$x$ .

The sum of all probabilities is 1:  $\sum_x P(X = x) = 1$

Example: Consider a fair six-sided die.

Let  $X$  be the random variable representing the result of a roll.

a) Possible values:  $X = \{1, 2, 3, 4, 5, 6\}$

b) PMF:  $P(X = x) = \frac{1}{6}$  for  $x = 1, 2, \dots, 6$ .

## Probability Density Function (PDF)

The PDF is used for continuous random variables. It describes the relative likelihood of the random variable taking on a particular value.

Properties:

1)  $f(x) \geq 0$  for all  $x$ .

2) The total area under the PDF curve is 1:  $\int_{-\infty}^{\infty} f(x) dx = 1$

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

3) Unlike the PMF, the PDF does not give the probability at a specific point but rather the probability within an interval:

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

$dx$

c) To find  $P(-1 \leq X \leq 1)$ , integrate the PDF over that range:

$$P(-1 \leq X \leq 1) = \int_{-1}^1 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$

$2\pi$

$$1 - e^{-2x^2}$$

dx

**Q2: What is Cumulative Density Function (CDF)? Explain with an example. Why CDF is used?**

The Cumulative Density Function (CDF) of a random variable  $X$  gives the probability that  $X$  will take a value less than or equal to a specific value  $x$

$x$ .

For both discrete and continuous random variables, the CDF is defined as:

$$F(x) = P(X \leq x)$$

## Properties of CDF

1)  $0 \leq F(x) \leq 1$  for all  $x$

$x$ .

2) The CDF is a non-decreasing function.

3) For a continuous random variable:  $F(x) = \int_{-\infty}^x f(t) dt$

where  $f(t)$

$f(t)$  is the Probability Density Function (PDF).

4) For a discrete random variable:  $F(x) = \sum_{t \leq x} P(X=t)$

## Example

## 1. Discrete Case: Rolling a Die

Let  $X$  represent the outcome of rolling a fair six-sided die. The PMF is:

$$P(X=x) = \frac{1}{6}, x \in \{1, 2, 3, 4, 5, 6\}$$

The CDF  $F(x) = P(X \leq x)$  is:

$$F(x) = \begin{cases} 0 & \text{if } x < 1 \\ \frac{1}{6} & \text{if } 1 \leq x < 2 \\ \frac{2}{6} & \text{if } 2 \leq x < 3 \\ \frac{3}{6} & \text{if } 3 \leq x < 4 \\ \frac{4}{6} & \text{if } 4 \leq x < 5 \\ \frac{5}{6} & \text{if } 5 \leq x < 6 \\ 1 & \text{if } x \geq 6 \end{cases}$$

0 6 1

6 2

6 3

6 4

6 5

1

if  $x < 1$  if  $1 \leq x < 2$  if  $2 \leq x < 3$  if  $3 \leq x < 4$  if  $4 \leq x < 5$  if  $5 \leq x < 6$  if  $x \geq 6$

## 2. Continuous Case: Normal Distribution

For a standard normal random variable  $Z$  with PDF:

$$f(z) = \frac{1}{\sqrt{2\pi}}$$

$$e^{-\frac{z^2}{2}}, -\infty < z < \infty$$

The CDF  $F(z)$  is:

$$F(z) = P(Z \leq z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt$$

$$-\infty$$

$$2\pi$$

$$e^{-\frac{t^2}{2}} dt$$

For example,  $F(0) = 0.5$   $F(0)=0.5$ , since half the area under the standard normal curve lies to the left of  $z = 0$   $z=0$ .

## Why is CDF Used?

1) Cumulative Probability: It allows us to calculate the probability of a random variable being less than or equal to a specific value, making it useful for probabilistic analyses.

2) Probability Intervals: We can use the CDF to calculate probabilities over intervals:  $P(a \leq X \leq b) = F(b) - F(a)$

3) Universal Representation: The CDF fully describes the distribution of a random variable and can be used to derive both PMF (for discrete variables) and PDF (for continuous variables).

4) Quantile Calculation: CDFs are used to find quantiles or percentiles of a distribution.

Q3: What are some examples of situations where the normal distribution might be used as a model? Explain how the parameters of the normal distribution relate to the shape of the distribution.

## Examples of Situations Where the Normal Distribution is Used

The normal distribution is widely used in various fields due to its prevalence in natural and human-made phenomena. Examples include:

- 1) Heights of Individuals: The heights of people in a population tend to follow a normal distribution, with most individuals clustering around the average height and fewer at the extremes.
- 2) IQ Scores: Intelligence quotient (IQ) scores are often modeled as a normal distribution with a mean of 100 and a standard deviation of 15.
- 3) Measurement Errors: Errors in scientific measurements or instruments typically follow a normal distribution because of the central limit theorem.
- 4) Stock Market Returns: Daily changes in stock prices or financial indices are approximately normally distributed in short periods.
- 5) Blood Pressure Levels: In a population, blood pressure measurements tend to exhibit a normal distribution.
- 6) Exam Scores: In a large class or standardized tests, scores often follow a normal distribution due to the aggregation of many small factors influencing performance.

## Parameters of the Normal Distribution and Their Effects on Shape

The normal distribution is characterized by two parameters: the mean ( $\mu$ ) and the standard deviation ( $\sigma$ ).

### 1) Mean ( $\mu$ )

- a) The mean determines the location of the center of the distribution.

b) It is the peak of the bell curve and represents the average value of the data.

c) Changing  $\mu$  shifts the curve left or right without altering its shape.

## 2) Standard Deviation ( $\sigma$ )

a) The standard deviation controls the spread or width of the distribution.

b) Larger  $\sigma$ : The curve becomes wider and flatter, indicating more variability in the data.

c) Smaller  $\sigma$ : The curve becomes narrower and taller, indicating less variability.

## Shape Effects

1) Symmetry: The normal distribution is symmetric about its mean  $\mu$ .

2) Bell Curve: The shape remains bell-like regardless of  $\mu$  or  $\sigma$ , but the specific characteristics of the bell depend on these parameters.

## Key Formula:

The probability density function (PDF) of the normal distribution is given by:

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

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

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

a)  $\mu$ : Controls the center of the distribution.

b)  $\sigma$ : Appears in both the denominator (scaling the height) and exponent (controlling spread), affecting the shape significantly.

## ## Visualizing the Effect of Parameters

### 1) Changing $\mu$ :

a)  $\mu=0$ : Symmetrical about 0.

b)  $\mu=5$ : Symmetrical about 5, shifted to the right.

### 2) Changing $\sigma$ :

a)  $\sigma=1$ : Narrow and peaked.

b)  $\sigma=3$ : Wider and flatter.

The normal distribution is used extensively because of its mathematical properties, its natural occurrence in many processes, and its flexibility in modeling data through  $\mu$  and  $\sigma$ .

## Q4: Explain the importance of Normal Distribution. Give a few real-life examples of Normal Distribution.

### #Importance of the Normal Distribution

The normal distribution is one of the most important probability distributions in statistics and is widely used due to its natural occurrence in numerous real-world scenarios. Its importance arises from the following:

#### 1. Central Limit Theorem (CLT):

1.1 The CLT states that the distribution of the sum or average of a large number of independent, identically distributed random variables tends toward a normal distribution, regardless of the original distribution of the data.

1.2 This property makes the normal distribution a cornerstone in statistical analysis and hypothesis testing.

#### 1. Simplification of Complex Problems:

2.1 Many real-world processes, when influenced by a large number of small, independent factors, result in outcomes that are approximately normally distributed. This allows for simplified modeling and analysis.

#### 1. Basis for Statistical Inference:

3.1 Statistical methods like confidence intervals, hypothesis testing, and regression analysis often assume a normal distribution of the data or residuals, making it foundational in inferential statistics.

#### 2. Widely Applicable Model:

4.1 The normal distribution is a good approximation for many types of data, making it versatile and easy to work with in fields like physics, finance, biology, and social sciences.

## Real-Life Examples of Normal Distribution

1. Human Heights: Heights of people within a specific population often follow a normal distribution, with most individuals clustering around the mean height.
2. IQ Scores: Intelligence quotient (IQ) scores are designed to follow a normal distribution with a mean of 100 and a standard deviation of 15. This ensures consistent assessment across populations.
3. Measurement Errors: Errors in scientific experiments and instruments, due to random variations, tend to follow a normal distribution.

4. **Stock Market Returns:**Daily returns of stocks or financial indices are approximately normally distributed in the short term, aiding risk management and portfolio analysis.
5. **Test Scores:**Scores on standardized exams like SATs or GREs are often normalized to follow a normal distribution, enabling fair comparison across test-takers.
6. **Blood Pressure Levels:**Blood pressure readings within a healthy population often form a bell-shaped curve, centered around the average normal blood pressure.
7. **Weight of Products:**In manufacturing, the weight of items like packaged food or bottled drinks tends to follow a normal distribution due to production process variability.

## Why Is the Normal Distribution Important in Practice?

- 1) **Predictability:**Many real-world variables that follow a normal distribution allow us to make predictions about future occurrences using probabilities.
- 2) **Standardization and Comparisons:**Data can be transformed into the standard normal distribution (z-scores) for easier comparison and interpretation.
- 3) **Decision-Making:**In business, healthcare, and policy-making, understanding whether data fits a normal distribution helps make informed decisions.
- 4) **Natural Occurrence:**The normal distribution naturally arises in processes influenced by numerous small, independent factors, making it a reliable model in diverse fields.

In summary, the normal distribution is central to understanding, modeling, and analyzing real-world data and phenomena, enabling robust and meaningful insights across a wide range of applications.

### **Q5: What is Bernoulli Distribution? Give an Example. What is the difference between Bernoulli Distribution and Binomial Distribution?**

**Bernoulli Distribution:** The Bernoulli distribution is a discrete probability distribution that describes the outcome of a single experiment with two possible outcomes: success (1) or failure (0). It is named after the Swiss mathematician Jacob Bernoulli.

#### **Probability Mass Function (PMF):**

$$P(X=x) = \begin{cases} p, & \text{if } x=1 \\ 1-p, & \text{if } x=0 \end{cases}$$

if  $x=1$  if  $x=0$

where  $0 \leq p \leq 1$ ,  $p$  is the probability of success, and  $1 - p$  is the probability of failure.

**Mean (Expected Value):**  $E[X] = p$

3) **Variance:**  $\text{Var}(X) = p(1-p)$   $\text{Var}(X)=p(1-p)$

**Example of Bernoulli Distribution:**

1) **Scenario: Flipping a coin.** a)  $X=1$  if the coin lands on heads. b)  $X=0$  if the coin lands on tails. c)  $p=0.5$  (assuming a fair coin). d) The outcome of a single coin flip follows a Bernoulli distribution.

## ***Difference Between Bernoulli and Binomial Distribution:***

Aspect	Bernoulli Distribution	Binomial Distribution
Definition	Describes the outcome of a single trial with two outcomes.	Describes the number of successes in $n$ independent Bernoulli trials.
Number of Trials	One trial ( $n = 1$ ).	Multiple trials ( $n \geq 1$ ).
PMF	$P(X = x) = p^x(1-p)^{1-x}$ .	$P(X = k) = \binom{n}{k}p^k(1-p)^{n-k}$ , where $k$ is the number of successes.
Random Variable	Binary ( $X = 0$ or $X = 1$ ).	Integer ( $X = 0, 1, 2, \dots, n$ ).
Mean	$p$ .	$n \cdot p$ .
Variance	$p(1-p)$ .	$n \cdot p \cdot (1-p)$ .
Example	Tossing a coin once.	Tossing a coin $n$ times and counting the heads.

**Example of Binomial Distribution:**

1) **Scenario:** Tossing a coin 5 times and counting the number of heads. - Each coin toss is a Bernoulli trial with  $p = 0.5$ . **The total number of heads follows a binomial distribution with  $n=5$  and  $p = 0.5$ \*\*.** Possible outcomes: 0,1,2,3,4,5 heads.

**Q6. Consider a dataset with a mean of 50 and a standard deviation of 10. If we assume that the dataset is normally distributed, what is the probability that a randomly selected observation will be greater than 60? Use the appropriate formula and show your calculations.**

To calculate the probability that a randomly selected observation from a normal distribution is greater than 60, we use the Z-score formula and the standard normal distribution table (or a computational tool). Here are the steps:

**Given: \*\***

- Mean ( $\mu$ ) = 50
- Standard Deviation ( $\sigma$ ) = 10

**Step 1: Calculate the Z-score**



The Z-score formula is:  $Z = (X - \mu) / \sigma$

Substitute the values:

$$Z = (60 - 50) / 10 = 10 / 10 = 1$$

### Step 2: Find the probability corresponding to $Z=1$

Using a standard normal distribution table or computational tools:

- The cumulative probability for  $Z=1$  is 0.8413. This is the probability that  $X \leq 60$ .

### Step 3: Calculate the probability for $X > 60$

The probability for  $X > 60$  is the complement of  $X \leq 60$ :  $P(X > 60) = 1 - P(X \leq 60)$

**Final Answer:** The probability that a randomly selected observation is greater than 60 is 0.1587 (or 15.87%).

## Q7: Explain uniform Distribution with an example.

### Uniform Distribution

The Uniform Distribution is a type of probability distribution in which all outcomes are equally likely. It can be described by two parameters:

**1) Lower Bound (a): The minimum value.**

**2) Upper Bound (b): The maximum value.**

For a uniform distribution, the probability density function (PDF) is constant across the range  $[a, b]$ , and zero outside this range.

### Formula for Uniform Distribution

If  $X$  is a random variable with a uniform distribution in the interval  $[a, b]$ , the PDF is:

$$f(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$$

The cumulative distribution function (CDF) is:

Properties

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

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

### Example

Suppose a train arrives at a station sometime between 2:00 PM and 2:30 PM. If the arrival time is equally likely throughout this 30-minute interval, the arrival time follows a uniform distribution. Parameters:

a=0 (minutes past 2:00 PM)  
b=30 (minutes past 2:00 PM)

PDF:  $f(x)=1/(b-a)=1/(30-0)=1/30, 0 \leq x \leq 30$

### Example:

**Case 1:** Discrete Uniform Distribution Suppose you roll a fair six-sided die. The possible outcomes are {1,2,3,4,5,6}{1,2,3,4,5,6}, and each has an equal probability of 1/6.

**Case 2:** Continuous Uniform Distribution Imagine a random number generator that picks a number between 2 and 8.

Here:

a) a=2, b=8

b) Probability of any specific subinterval (e.g., 3 to 5) is proportional to its length.

If we want to find the probability of selecting a number between 4 and 6:

$$**P(4 \leq X \leq 6)** = (6-4) \cdot 1/(b-a) = 2 \cdot (1/6) = (1/3) .$$

### Key Properties:

- The graph of the PDF for a continuous uniform distribution is a horizontal line.
- For a discrete case, the PMF is constant for all outcomes.

Uniform distribution is commonly used in simulations and scenarios where all outcomes are equally plausible within a range.

### Q-8: What is the z score? State the importance of the z score ?

The z-score (or standard score) is a statistical measure that describes how far a data point is from the mean of a dataset, measured in terms of standard deviations. It indicates whether a data point is above or below the mean and by how many standard deviations.

The formula for the z-score is:

$$**z=(x-\mu)/\sigma **$$

Where:

x = the value of the data point  
 $\mu$  = the mean of the dataset  
 $\sigma$  = the standard deviation of the dataset

A Z-score (or standard score) measures how many standard deviations a data point is from the mean of a dataset.

### Importance of Z-Score:

- 1. Standardization:** Z-scores transform data from different scales to a common scale, making it easier to compare different datasets.
- 2. Outlier Detection:** Data points with high absolute Z-scores (e.g., greater than 3 or less than -3) are typically considered outliers.
- 3. Probability Calculations:** Z-scores are essential in statistics for calculating probabilities and areas under the normal distribution curve.
- 4. Statistical Testing:** They play a crucial role in hypothesis testing, such as in z-tests, to assess whether a sample mean differs significantly from a known population mean.
- 5. Data Normalization:** Useful in machine learning models to normalize features for better performance during training and prediction.
- 6. Risk Management:** In finance, Z-scores help assess the likelihood of extreme price movements or financial distress.
- 7. Quality Control:** Z-scores help monitor and control processes by identifying variations in production or operational systems.

In essence, the Z-score is a versatile tool for understanding and interpreting the distribution and behavior of data in various fields.

### Q9: What is Central Limit Theorem? State the significance of the Central Limit Theorem.

#### Central Limit Theorem (CLT)

The **Central Limit Theorem (CLT)** states that when independent random samples of sufficiently large size ( $n$ ) are drawn from any population with a finite mean ( $\mu$ ) and variance ( $\sigma^2$ ), the sampling distribution of the sample mean will approximate a normal distribution, regardless of the shape of the original population distribution.

#### Mathematical Expression:

If  $X_1, X_2, \dots, X_n$  are independent and identically distributed random variables with mean  $\mu$  and standard deviation  $\sigma$ , then for large  $n$ , the distribution of the sample mean  $\bar{X}$  approaches:

$$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

#### Significance of the Central Limit Theorem:

- 1) Normal Approximation:** The CLT allows us to use normal distribution techniques (e.g., z-scores) to make inferences about sample means, even if the original population is not normally distributed.
- 2) Foundation of Inferential Statistics:** It underpins many statistical methods, such as confidence intervals and hypothesis testing.

**3) Sampling Efficiency:**As sample size increases, the distribution of the sample mean becomes more normal, making it easier to predict and analyze.

**4) Real-World Applications:**Used in various fields, including finance (stock price analysis), manufacturing (quality control), and research (survey analysis).

**5) Reduced Data Dependency:**It shows that large samples provide reliable estimates of population parameters, making statistical methods applicable across diverse datasets.

#### **Conditions for CLT Application:**

- The sample size  $n$  should typically be  $n \geq 30$  for the theorem to hold well.
- Samples must be independent.
- The population from which the sample is drawn must have finite variance.

The CLT is a powerful concept that simplifies complex real-world problems by enabling the application of well-understood statistical techniques.

#### **Q10: State the assumptions of the Central Limit Theorem.**

##### **Assumptions of the Central Limit Theorem (CLT):**

**1. Independence of Samples:**The samples must be independent of each other. The value of one observation should not influence another.

**2. Identically Distributed Random Variables:**The data points should come from the same distribution and be identically distributed if drawn from a finite population.

**3. Finite Variance:**The population from which samples are drawn should have a finite variance ( $\sigma^2 < \infty$ ) to ensure the sample mean's distribution stabilizes.

**4. Sample Size:**A sufficiently large sample size ( $n$ ) is necessary. Typically,  $n \geq 30$  is considered large enough, but for skewed distributions or those with heavy tails, a larger sample may be required.

**5. Random Sampling:**The samples should be drawn randomly to ensure that they accurately represent the population.

##### **Special Considerations:**

- When sampling from a normal distribution, the CLT holds regardless of sample size.
- For non-normal populations, larger sample sizes improve the approximation of the sample mean to a normal distribution.
- If the population distribution is heavily skewed or has outliers, stricter adherence to sample size guidelines is required for the theorem to hold.

These assumptions ensure the robustness of the CLT and its broad application in statistical methods.

