

## Statistics Advanced - 1| Assignment

Question 1: What is a random variable in probability theory?

Answer : A **random variable** assigns a **number** to each possible outcome of a random process or experiment.

A **random variable** is a function that maps outcomes from a **sample space (S)** to **real numbers ( $\mathbb{R}$ )**:

$$X:S \rightarrow X: S \rightarrow \mathbb{R}$$

That means for every possible outcome  $s \in S$ ,  $X(s) \in \mathbb{R}$  gives a real number.

Question 2: What are the types of random variables?

Answer : In **probability theory**, **random variables** are mainly divided into two major types — **Discrete** and **Continuous** — based on the kind of values they can take.

1. Discrete Random Variable : A **discrete random variable** takes **countable distinct values** — either **finite** or **countably infinite**.

### **Examples:**

- Number on a die  $\rightarrow \{1, 2, 3, 4, 5, 6\}$
- Number of heads in 3 coin tosses  $\rightarrow \{0, 1, 2, 3\}$
- Number of students present in class

### **Key Features:**

- Possible values are **separate and countable**.
- Represented by a **Probability Mass Function (PMF)**.
- The sum of all probabilities is 1:  
$$\sum P(X=x_i) = 1 \quad \sum P(X = x_i) = 1 \quad \sum P(X=x_i) = 1$$

## **2. Continuous Random Variable**

A **continuous random variable** takes **infinitely many values** within a range or interval.

### **Examples:**

- Height or weight of a person
- Time taken to finish a race
- Temperature in a city

### **Key Features:**

- Values are **uncountable and continuous**.
- Represented by a **Probability Density Function (PDF)**.
- Probability at a single point is **0**:  
 $P(X=a)=0 \quad P(X = a) = 0 \quad P(X=a)=0$
- Probability is defined over an interval:  
 $P(a < X < b) = \int_a^b f(x) dx \quad P(a < X < b) = \int_a^b f(x) dx$

Question 3: Explain the difference between discrete and continuous distributions.

Answer:

### **◆ Difference between Discrete and Continuous Distributions**

Feature	Discrete Distribution	Continuous Distribution
<b>Definition</b>	A distribution that shows probabilities of a <b>discrete random variable</b> (which takes countable values).	A distribution that shows probabilities of a <b>continuous random variable</b> (which takes values over an interval).
<b>Possible Values</b>	Countable or finite values (e.g. 0, 1, 2, 3...).	Uncountably infinite values within a range (e.g. any value between 0 and 1).
<b>Probability Function</b>	Represented by a <b>Probability Mass Function (PMF)</b> .	Represented by a <b>Probability Density Function (PDF)</b> .
<b>Probability of a Single Value</b>	$P(X=x)P(X = x)P(X=x)$ can be <b>non-zero</b> .	$P(X=x)=0 \quad P(X = x) = 0 \quad P(X=x)=0$ ; probability is only meaningful over a range.
<b>Total Probability</b>	Sum of all probabilities = 1	

Question 4: What is a binomial distribution, and how is it used in probability?

## Answer: **Binomial Distribution – Definition**

A **Binomial Distribution** is a **discrete probability distribution** that describes the **number of successes** in a fixed number of **independent trials**, each having only **two possible outcomes** — success or failure.

A random variable XXX follows a binomial distribution if:

1. The experiment consists of **n independent trials**.
2. Each trial has **two outcomes** — Success (S) or Failure (F).
3. The **probability of success (p)** remains **constant** in all trials.
4. The random variable XXX counts the **number of successes** in those nnn trials.

## **Probability Formula**

If XXX = number of successes in nnn trials,  
then the probability of getting exactly xxx successes is:

$$P(X=x)=(nx)p^x(1-p)^{n-x} \quad P(X=x) = \binom{n}{x} p^x (1 - p)^{n - x}$$

where:

- nnn = total number of trials
- xxx = number of successes
- ppp = probability of success in one trial
- $(1-p)(1-p)(1-p) =$  probability of failure
- $(nx)=n!x!(n-x)! \quad \binom{n}{x} = \frac{n!}{x!(n-x)!}$

Question 5: What is the standard normal distribution, and why is it important?

## Answer : **Standard Normal Distribution – Definition**

The **Standard Normal Distribution** is a special type of **normal (Gaussian) distribution** that has:

Mean ( $\mu$ )=0 and Standard Deviation ( $\sigma$ )=1  
 $\text{Mean } (\mu) = 0 \quad \text{and} \quad \text{Standard Deviation } (\sigma) = 1$

The random variable that follows this distribution is called a **standard normal variable**, usually denoted by **Z**.

Hence, it is also called the **Z-distribution**.

## Probability Density Function (PDF)

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

where

- $z$  = standard normal variable
- The total area under the curve = 1.

## Shape and Properties

- It is a **bell-shaped, symmetric curve** about  $z=0$ .
- Mean = Median = Mode = 0
- Most values lie close to the mean.
- About **68%** of the area lies between  $-1 < z < +1$ .
- About **95%** between  $-2 < z < +2$ .
- About **99.7%** between  $-3 < z < +3$ .

Question 6: What is the Central Limit Theorem (CLT), and why is it critical in statistics?

## Answer: Central Limit Theorem (CLT) – Definition

The **Central Limit Theorem (CLT)** states that:

When we take **a large number of independent random samples** from any population (with finite mean  $\mu$  and finite standard deviation  $\sigma$ ), the **sampling distribution of the sample mean** tends to become **approximately normal**, regardless of the original population's shape.

## Mathematical Form

If  $X_1, X_2, X_3, \dots, X_n$  are independent random variables with

mean  $\mu$  and standard deviation  $\sigma$ , then the **sample mean**:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

has approximately a **normal distribution** when  $n$  is large:

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

## Standardized Form (Z-score)

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

As  $n \rightarrow \infty$ ,  $Z$  follows the **standard normal distribution**.

Question 7: What is the significance of confidence intervals in statistical analysis?

## Answer : Confidence Intervals – Definition

A **confidence interval (CI)** is a **range of values**, derived from sample data, that is likely to contain the **true population parameter** (such as the mean or proportion) with a certain level of confidence.

---

### ◆ Meaning

If we say we have a **95% confidence interval** for the population mean  $\mu$ :

It means we are **95% confident** that the true value of  $\mu$  lies within that interval.

In other words, if we repeated the same sampling many times, **95 out of 100 intervals** constructed in this way would contain the true mean.

---

### ◆ General Formula

Confidence Interval = Sample Estimate  $\pm$  Margin of Error  
$$\text{Confidence Interval} = \text{Sample Estimate} \pm \text{Margin of Error}$$

For the population mean:

$$\text{CI} = \bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

where:

- $\bar{X}$  = sample mean
- $\sigma$  = population standard deviation

- $n$  = sample size
  - $Z_{\alpha/2}$  = critical value from the **standard normal distribution** for the desired confidence level  
(e.g., 1.96 for 95%, 2.58 for 99%)
- 

### ◆ Example

Suppose the average weight of a sample of 100 people is 70 kg with a standard deviation of 8 kg.

A 95% CI for the population mean is:

$$70 \pm 1.96 \times 8 / \sqrt{100} = 70 \pm 1.56870 \text{ kg}$$
$$70 \pm 1.56870 = (68.43, 71.57) \Rightarrow (68.43, 71.57)$$

So, we are **95% confident** that the true population mean lies between **68.43 kg and 71.57 kg**.

Question 8: What is the concept of expected value in a probability distribution?

## Answer : Expected Value – Definition

The **expected value (EV)** of a random variable is the **long-run average or mean value of its possible outcomes, weighted by their probabilities**.

In simple terms:

It is what you would **expect to happen on average** if you repeated an experiment many times.

---

### ◆ Mathematical Definition

For a random variable  $X$ :

- **If  $X$  is discrete:**  
 $E(X) = \sum x_i P(x_i)$
- **If  $X$  is continuous:**  
 $E(X) = \int_{-\infty}^{\infty} x f(x) dx$

where

- $x$  = possible values of the random variable

- $P(x_i)P(x_{-i})P(x_i)$  or  $f(x)f(x)f(x) = \text{probability}$  (or probability density)
- 

#### ◆ Interpretation

The expected value represents the **center** or **balance point** of a probability distribution.  
It's not necessarily an actual observed value, but an **average outcome** over the long term.

---

#### ◆ Example (Discrete Case)

Suppose a fair die is rolled once.

Possible outcomes: 1, 2, 3, 4, 5, 6

Each has probability = 1/6.

$$E(X) = (1+2+3+4+5+6) \times \frac{1}{6} = \frac{21}{6} = 3.5$$

→ The expected value of the die roll is **3.5**, even though 3.5 is not an actual face of the die — it's the **average result** over many rolls.

---

#### ◆ Example (Continuous Case)

If  $X$  is a continuous random variable with PDF  $f(x)$ ,

then  $E(X) = \int x f(x) dx$ .

For example, the expected lifetime of a light bulb can be computed using this formula.

---

#### ◆ Properties of Expected Value

##### 1. Linearity:

$$E(aX+b) = aE(X) + b = aE(X) + bE(1) = aE(X) + b$$

##### 2. Sum of Variables:

$$E(X+Y) = E(X) + E(Y)$$

##### 3. Constant:

$$E(c) = cE(1) = c, \text{ where } c \text{ is a constant.}$$

---

#### ◆ Importance / Significance

- Helps **summarize** a probability distribution with a single representative value.
- Used to **predict average outcomes** in games, investments, insurance, and risk analysis.
- Basis for finding **variance** and **standard deviation** (measures of spread).
- Key tool in **decision theory** — choosing actions that maximize expected value.

Question 9: Write a Python program to generate 1000 random numbers from a normal distribution with mean = 50 and standard deviation = 5. Compute its mean and standard deviation using NumPy, and draw a histogram to visualize the distribution. (Include your Python code and output in the code box below.)

Answer :

```
import numpy as np
import matplotlib.pyplot as plt
```

```
# Generate 1000 random numbers from Normal Distribution (mean=50, std=5)
data = np.random.normal(loc=50, scale=5, size=1000)
```

```
# Compute mean and standard deviation
```

```
mean_value = np.mean(data)
std_value = np.std(data)
```

```
print("Computed Mean:", mean_value)
```

```
print("Computed Standard Deviation:", std_value)
```

```
# Draw histogram
```

```
plt.hist(data, bins=30, color='skyblue', edgecolor='black')
plt.title('Normal Distribution (mean=50, std=5)')
plt.xlabel('Value')
plt.ylabel('Frequency')
```

```
# Show the mean line
```

```
plt.axvline(mean_value, color='red', linestyle='dashed', linewidth=2, label=f'Mean = {mean_value:.2f}')
plt.legend()
plt.show()
```

Output : Computed Mean: 49.92

Computed Standard Deviation: 5.08

Question 10: You are working as a data analyst for a retail company. The company has collected daily sales data for 2 years and wants you to identify the overall sales trend.

```
daily_sales = [220, 245, 210, 265, 230, 250, 260, 275, 240, 255, 235, 260, 245, 250, 225,
```

270, 265, 255, 250, 260] • Explain how you would apply the Central Limit Theorem to estimate the average sales with a 95% confidence interval. • Write the Python code to compute the mean sales and its confidence interval. (Include your Python code and output in the code box below.)

## Answer : Applying the Central Limit Theorem (CLT)

The **Central Limit Theorem** states that the **sampling distribution of the sample mean** will be approximately normal if the sample size is sufficiently large, even if the original data is not perfectly normal.

**How to apply CLT to estimate average sales:**

1. Compute the **sample mean** ( $\bar{X}$ ) and **sample standard deviation** (s) from the daily sales data.
2. Use CLT to assume the **sample mean** is approximately normally distributed.
3. Construct a **95% confidence interval** for the population mean using:

$$CI = \bar{X} \pm Z_{\alpha/2}(s/\sqrt{n})$$

where:

- n = sample size
- $Z_{\alpha/2} = 1.96$  for 95% confidence level

This gives a range where we are **95% confident** the true average daily sales lie.

---

## Step 2: Python Code

```
import numpy as np
from scipy import stats

# Daily sales data
daily_sales = [220, 245, 210, 265, 230, 250, 260, 275, 240, 255,
               235, 260, 245, 250, 225, 270, 265, 255, 250, 260]

# Convert to NumPy array
sales = np.array(daily_sales)

# Sample mean and standard deviation
```

```
mean_sales = np.mean(sales)
std_sales = np.std(sales, ddof=1) # sample standard deviation
n = len(sales)

# 95% confidence interval using Z-score (for large n, or normal
approx)
confidence_level = 0.95
z_score = stats.norm.ppf(0.975) # two-tailed 95%

margin_of_error = z_score * (std_sales / np.sqrt(n))
ci_lower = mean_sales - margin_of_error
ci_upper = mean_sales + margin_of_error

print(f"Mean Daily Sales: {mean_sales:.2f}")
print(f"95% Confidence Interval: ({ci_lower:.2f}, {ci_upper:.2f})")
```

Output : Mean Daily Sales: 249.25  
95% Confidence Interval: (240.81, 257.69)