

This is a assignment of Statistics Advance Part 1

1. What is a random variable in probability theory?

A **random variable** is a variable that takes **numerical values** based on the outcome of a **random experiment**.

- It assigns a number to each possible outcome.
- It helps us **quantify uncertainty** in probability.

Example:

Tossing a coin:

- If Head → assign 1
- If Tail → assign 0
Then the random variable X = outcome of coin toss:
- Possible values: 0 or 1

So, X is a **random variable** which can take different values depending on chance.

2. What are the types of random variables?

There are **two main types**:

A. Discrete Random Variable

- Takes **countable** values (finite or infinite but countable)
- Often results from **counting** something

Examples:

- Number of students in a class
- Number of heads in 5 coin tosses: {0, 1, 2, 3, 4, 5}
- Rolling a die: {1, 2, 3, 4, 5, 6}

Graph:

Bar graph with spikes at each value

B. Continuous Random Variable

- Takes **infinite values in a range**
- Often results from **measuring** something

Examples:

- Height of a person (like 160.5 cm, 161.2 cm, etc.)
- Time taken to complete a task
- Temperature readings

Graph:

Smooth **curve** (like the bell curve for normal distribution)

3. What is the difference between discrete and continuous distributions?

Discrete Distribution:

- Represents probabilities of **discrete random variables** (countable outcomes).
- Each individual value has a **non-zero probability**.
- The **sum** of all probabilities = 1

Example:

Rolling a die → outcomes: 1, 2, 3, 4, 5, 6

Each has probability = $1/6$

Graph:

Bar graph with **spikes** at each value.

Continuous Distribution:

- Represents probabilities of **continuous random variables** (infinite outcomes in a range).
- The probability of any exact value is **0**
- Probability is calculated over an **interval** (like between 2 and 4)
- The **area under the curve** = 1

Example:

Height of people → measured as 160.2 cm, 160.25 cm, etc. (infinite possibilities)

Graph:

A **smooth curve** (e.g., normal distribution bell curve)

Key Differences Table:

Feature	Discrete Distribution	Continuous Distribution
Random Variable Type	Discrete (countable)	Continuous (infinite values)
Probability of single value	> 0 (e.g., $P(X=2) = 0.2$)	0 (e.g., $P(X=2) = 0$)
Total Probability	Sum of all points = 1	Area under curve = 1
Example	Coin toss, dice roll	Height, weight, time
Visualization	Bar chart	Smooth curve (PDF)

4. What are probability distribution functions (PDF)?

A **Probability Distribution Function (PDF)** is a function that **describes the likelihood** of a random variable taking on a particular value.

There are **two types** based on the type of random variable:

A. PDF for Discrete Random Variables

It gives **$P(X = x)$** for every possible outcome x .

This is also called a **Probability Mass Function (PMF)**.

Example:

For a fair die:

$$P(X = 3) = 1/6$$

B. PDF for Continuous Random Variables

- It's a **curve** that represents **probability density**, not exact probabilities.
- The actual probability is calculated over an **interval** using **area under the curve**.

Example:

If X is a normal random variable,

$$P(2 \leq X \leq 3) = \text{Area under PDF curve from } x = 2 \text{ to } x = 3$$

For continuous variables, $P(X = a) = 0$
(because one point has zero width \Rightarrow no area under curve)

Properties of a PDF:

1. $\text{PDF}(x) \geq 0$ for all x
2. Total area under the curve = 1
3. For continuous case:
$$P(a \leq X \leq b) = \int_a^b \text{PDF}(x) dx$$

5. How do cumulative distribution functions (CDF) differ from probability distribution functions (PDF)?

Probability Distribution Function (PDF):

- **PDF** shows **how likely** a random variable is to take a specific value (or range of values).
- For **discrete variables**, it's called **PMF** (Probability Mass Function).
- For **continuous variables**, the PDF is a **curve** — the probability over an interval is the **area under the curve**.

Key Point:

- PDF gives **instantaneous probability density**.
 - In continuous case: probability at a single point is **0**.
 - Mathematically, for continuous random variable
 - $XP(a \leq X \leq b) = \text{Area under PDF from } a \text{ to } b$
-

Cumulative Distribution Function (CDF):

- **CDF** gives the **cumulative probability** that the random variable is **less than or equal to** a value.
- It **adds up** the probability from the lowest value up to **x**.

Key Point:

- CDF is **always increasing** from 0 to 1.
 - Mathematically:
 $F(x) = P(X \leq x)$
-

Difference Summary Table:

Feature	PDF	CDF
Meaning	Probability density at a value	Probability up to a value
Range of Values	Can be > 1 (for densities)	Always between 0 and 1
Shape	Bell curve, spikes, etc.	Always non-decreasing curve
Usage	Find density at point	Find cumulative probability
Relation	CDF is the integral of PDF	Derivative of CDF gives PDF

Simple Example (Dice Roll):

- PDF/PMF:
 - $P(X = 3) = 1/6$
 - CDF:
 - $P(X \leq 3) = P(X=1) + P(X=2) + P(X=3) = 1/6 + 1/6 + 1/6 = 0.5$
-

6. What is a discrete uniform distribution?

Discrete Uniform Distribution:

- A **discrete uniform distribution** is when **all possible outcomes** have the **same probability**.
 - Each outcome is **equally likely**.
-

Examples:

- Rolling a **fair six-sided die**:
Each side {1, 2, 3, 4, 5, 6} has probability = $1/6$
 - Picking a random day of the week:
{Monday, Tuesday, ..., Sunday} → Probability = $1/7$ each
-

Key Properties:

- **Equal probability** for all values
 - If there are n possible outcomes:
 $P(X=x_i)=1/n$
 - Mean (Expected Value):
 $E(X)=(a+b)/2$
 - Variance:
 $Var(X)=(b-a+1)^2-1/12$
where a and b are the minimum and maximum values.
-

Summary:

Feature	Discrete Uniform Distribution
Probability of each outcome	Equal ($1/n$)
Examples	Dice roll, random lottery numbers
Shape of PMF	Flat horizontal line (equal spikes)

7. What are the key properties of a Bernoulli distribution?

Bernoulli Distribution:

A **Bernoulli distribution** models a random experiment with only **two possible outcomes**:

- **Success** (usually coded as 1)
- **Failure** (usually coded as 0)

It is the **simplest discrete probability distribution**.

Examples:

- Tossing a coin → Head (1), Tail (0)
- Passing or failing a test
- Yes or No survey responses

Probability Function:

Let p = probability of success (1), then:

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

Key Properties:

Property		Value / Formula
Possible values		0 or 1
Mean (Expected value)		$E(X)=p$
Variance		$\text{Var}(X)=p(1-p)$
Skewness		$1-2p/\sqrt{p(1-p)}$
Distribution type		Discrete

Graph:

A bar graph with spikes at:

- $0 \rightarrow \text{height} = 1 - p$
- $1 \rightarrow \text{height} = p$

 **Use:**

Bernoulli is used as the **building block** for:

- **Binomial distribution**
 - **Logistic regression**
 - **Binary classification in ML**
-

8. What is the Binomial Distribution, and how is it used in probability?

Binomial Distribution:

The **Binomial Distribution** models the **number of successes** in a **fixed number of independent Bernoulli trials**, each with the same probability of success.

Real-Life Examples:

- Number of heads in 10 coin tosses
 - Number of correct answers in a multiple-choice test
 - Number of defective items in a batch
-

Key Parameters:

- n = number of trials
 - p = probability of success in a single trial
-

Probability Mass Function (PMF):

$$P(X=k) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$$

Where:

$\binom{n}{k}$ = combinations of k successes out of n trials

- k = number of successes ($0 \leq k \leq n$)

Key Properties:

Property	Formula / Value
Mean (Expected Value)	$E(X) = np$
Variance	$Var(X) = np(1-p)$
Distribution Type	Discrete
Shape	Bell-shaped (when n is large)

Conditions to Use Binomial:

1. Fixed number of trials n
 2. Each trial is **independent**
 3. Each trial has only **two outcomes**: success/failure
 4. Constant probability p across trials
-

Use in Probability:

- To find probability of **exact number of successes** (e.g., exactly 3 heads)
- To **model real-world situations** like quality control, marketing, medicine

- **Hypothesis testing** (binomial test)
-

Example:

Q: What's the probability of getting **2 heads** in **3 coin tosses**?

Here:

- $n = 3, p = 0.5, k = 2$
- $P(X=2) = \binom{3}{2} \cdot (0.5)^2 \cdot (0.5)^1 = 3 \cdot 0.25 \cdot 0.5 = 0.375$

9. What is the Poisson Distribution and Where is it Applied?

Poisson Distribution:

The **Poisson distribution** models the **number of times an event occurs** in a **fixed interval of time or space**, if:

- Events occur **independently**,
 - The **rate of occurrence is constant**, and
 - Two events cannot happen at **exactly the same instant**.
-

Formula (PMF):

If λ is the **average rate** (mean number of occurrences), and k is the **actual number of occurrences**, then:

$$P(X=k) = e^{-\lambda} \cdot \lambda^k / k!$$

Where:

- $e \approx 2.718$
- $\lambda > 0$
- $k = 0, 1, 2, \dots$

Key Properties:

Property	Value / Formula
Mean (Expected value)	λ
Variance	λ
Skewness	$1/\sqrt{\lambda}$
Distribution Type	Discrete
Support	$k \in \{0, 1, 2, \dots\}$

Applications:

- Number of calls at a call center per minute
 - Number of emails received per hour
 - Number of accidents at an intersection per day
 - Number of typing errors per page
 - Number of customers entering a shop in an hour
-

When to Use Poisson:

Use when you're counting **how many times** something happens in a **fixed interval** of:

- Time
- Space
- Area
- Volume

10. What is a Continuous Uniform Distribution?

Continuous Uniform Distribution:

A **continuous uniform distribution** is where **every value within a given interval** is **equally likely** to occur.

Example:

- Random number between 0 and 1
 - Choosing a time uniformly between 2 PM and 4 PM
-

PDF Formula:

If a variable X is uniformly distributed between a and b :

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

Key Properties:

Property

Formula

Range	$[a, b]$
Mean	$E(X) = a + b/2$
Variance	$Var(X) = (b-a)^2 / 12$
Probability	All values in $[a, b]$ equally likely
Distribution Type	Continuous

Applications:

- Simulating random wait times
 - Random sampling in simulations
 - Choosing a point at random on a straight line
 - Any situation where every outcome in a range is equally likely
-

Comparison with Discrete Uniform:

Feature	Discrete Uniform	Continuous Uniform
Values	Countable (e.g. dice)	Infinite (e.g. time)

Probability Specific values Intervals (area under PDF)

Graph Bar chart Flat horizontal line

11. What are the Characteristics of a Normal Distribution?

Normal Distribution:

The **normal distribution** is a **bell-shaped** and **symmetric** probability distribution used to describe many **natural phenomena**.

Key Characteristics:

Feature	Description
Shape	Bell-shaped curve (Gaussian)
Symmetry	Perfectly symmetric about the mean
Mean = Median = Mode	All three are equal and lie at the center
Tails	Extend infinitely in both directions but never touch the x-axis
Spread	Controlled by standard deviation (σ)

68-95-99.7 Rule

Describes how data falls within standard deviations from the mean

68-95-99.7 Rule:

- 68% of data lies within $\pm 1\sigma$
- 95% of data lies within $\pm 2\sigma$
- 99.7% of data lies within $\pm 3\sigma$

Mathematical Formula (PDF):

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

Where:

- μ = mean
 - σ = standard deviation
-

Real-World Examples:

- Heights, weights, IQ scores, exam scores, etc.
-

12. What is the Standard Normal Distribution, and Why is It Important?

Standard Normal Distribution:

A **standard normal distribution** is a **special case** of the normal distribution with:

- **Mean (μ) = 0**
 - **Standard deviation (σ) = 1**
-

Importance:

1. **Used for Z-scores:**

Transforms any normal variable into a standard form:

$$Z = \frac{X - \mu}{\sigma}$$

- 2.

This allows comparisons between **different normal distributions**.

3. **Simplifies calculations:**

Most statistical tables (like probability charts) are based on the standard normal distribution.

4. **Foundation for inferential statistics:**

It's the core of **hypothesis testing**, **confidence intervals**, and **p-values**.

Visual:

The graph is bell-shaped and centered at **0**, with spread controlled by standard deviation = 1.

13. What is the Central Limit Theorem (CLT), and Why is It Critical in Statistics?

Central Limit Theorem (CLT):

The **CLT states** that:

When you take **many random samples** from any population (with any shape) and calculate the **mean** of each sample, the **distribution of those sample means** will tend to be **normal**, if the **sample size is large enough**.

Key Points:

- The **original data** can be skewed or non-normal.
- But the **sampling distribution of the mean** becomes **normal** as n increases.
- Usually, $n \geq 30$ is considered "large enough."

Why is it important?

It **justifies using normal-based methods** (like Z-tests, t-tests) even if the data isn't perfectly normal.

Forms the **foundation of inferential statistics**: estimates, confidence intervals, and hypothesis tests.

Formula (Sampling Distribution):

Mean of sample means = μ

Standard deviation (standard error) = σ / \sqrt{n}

14. How Does the Central Limit Theorem Relate to the Normal Distribution?

Connection:

- The **CLT explains why** the **normal distribution appears so often** in statistics.
- It shows that the **sampling distribution of the mean** will be **normal**, even if the original data is not.
- It gives the **normal distribution** a **central role** in:
 - Hypothesis testing
 - Confidence intervals
 - Statistical modeling

-

Visual Summary:

- | | | |
|-----------------|---------------|-----------------------|
| ○ Original data | sample size | sampling distribution |
| ○ Skewed | Small (n=5) | not normal |
| ○ Skewed | Large (n=30+) | Approaches Normal |
| ○ Normal | Any n | Always Normal |

15. What is the Application of Z-Statistics in Hypothesis Testing?

Z-statistic (Z-test):

The Z-statistic is used in hypothesis testing when:

- The population standard deviation (σ) is known.
- The sample size is large ($n \geq 30$).

Applications of Z-Statistic:

1. One-sample Z-test:

- Compare the sample mean to a known population mean.

2. Two-sample Z-test:

- Compare the means of two independent large samples.

3. Z-test for proportions:

- Compare sample proportion to a known population proportion.

Z-Test Process:

1. State null (H_0) and alternative (H_1) hypotheses.
 2. Calculate the Z-score.
 3. Find the critical value or p-value.
 4. Make a decision:
 - If $|Z| > \text{critical value} \rightarrow \text{reject } H_0$.
 - If $\text{p-value} < \alpha$ (e.g. 0.05) $\rightarrow \text{reject } H_0$.
-

Example Use Cases:

- Quality control
 - Medical trials
 - Comparing exam performance to national average
-

16. How Do You Calculate a Z-Score, and What Does It Represent?

Z-score:

The Z-score tells you how many standard deviations a value is from the mean.

Formula: $Z = \frac{X - \mu}{\sigma}$

Where:

- X = individual data point
- μ = mean
- σ = standard deviation

Interpretation:

Z-Score	Meaning
0	Exactly at the mean
+1	1 standard deviation above the mean
-1	1 standard deviation below the mean
± 2 or more	Far from the mean (outliers)

Use in Statistics:

- Comparing scores across different distributions
- Outlier detection
- Hypothesis testing (Z-tests)
- Normalization in machine learning

17. What Are Point Estimates and Interval Estimates in Statistics?

Point Estimate:

A single value used to estimate a population parameter.

Example:

Sample mean \bar{x} is a point estimate of population mean μ

Interval Estimate:

A range of values used to estimate the population parameter, along with a confidence level.

Example:

"We are 95% confident the population mean lies between 68 and 72."

Comparison:

Point estimate single best guess $\bar{x}=70$

Interval estimate range with confidence $68 \leq \mu \leq 72$ with 95% confidence

18. What is the Significance of Confidence Intervals in Statistical Analysis?

Confidence Interval (CI):

A range around a point estimate that likely contains the true population parameter with a specified confidence level (e.g., 95%).

Why Important?

1. **Provides uncertainty range:**
Shows how much variability is in your estimate.
2. **More informative than point estimate:**
Tells you *how precise* your estimate is.

3. **Used in decision-making:**
Helps determine statistical significance in hypothesis testing.
 4. **Linked to hypothesis testing:**
If the hypothesized value (like μ_0) lies outside the confidence interval, we reject the null hypothesis.
-

Common Confidence Levels:

- 90%
 - 95% (most common)
 - 99%
-

Formula (for mean, known σ):

$$CI = \bar{x} \pm Z_{\alpha/2} \cdot \sigma / \sqrt{n}$$

19. What is the Relationship Between a Z-Score and a Confidence Interval?

Z-Score and Confidence Interval (CI) are closely connected.

- A **Z-score** determines **how far** a sample mean is from the population mean in terms of **standard deviations**.
- A **Confidence Interval** uses the **Z-score** as a **critical value** to define the range where the true parameter is likely to lie.

Formula (for known σ):

$$CI = \bar{x} \pm Z_{\alpha/2} \cdot \sigma / \sqrt{n}$$

\bar{x} = sample mean

$Z_{\alpha/2}$ = Z-score corresponding to the desired confidence level (e.g. 1.96 for 95%)

= standard error σ/\sqrt{n}

The **Z-score** sets the width of the confidence interval. Higher confidence level → higher Z-score → wider interval.

20. How Are Z-Scores Used to Compare Different Distributions?

Z-score standardizes values from different distributions to a common scale.

$$Z = \frac{X - \mu}{\sigma}$$

- This lets you compare values **even if** the means and standard deviations of the distributions are **different**.

Example:

- Alice scores 85 on a test with mean 80 and SD 5 → $Z = 1$
- Bob scores 90 on a test with mean 88 and SD 4 → $Z = 0.5$

Alice performed **better relative to her group** ($Z=1$ vs $Z=0.5$).

Use Cases:

- Comparing test scores across different subjects
 - Outlier detection
 - Normalization in machine learning
-

21. What Are the Assumptions for Applying the Central Limit Theorem (CLT)?

The **Central Limit Theorem (CLT)** allows us to assume that the **sampling distribution of the sample mean** is approximately **normal**, **regardless of population shape**, if certain assumptions are met.

- **Key Assumptions:**

- | Assumption | Description |
|-------------------------------|--|
| • Random Sampling | Data must be selected randomly and independently. |
| • Sample Size ($n \geq 30$) | A large enough sample size is required ($n \geq 30$ is a rule of thumb). Smaller n is okay if the population is already normal. |
| • Finite Variance | The population should have a finite standard deviation (σ). |
| • | |

Conclusion:

If these assumptions are satisfied, we can use **normal distribution methods** (like Z-tests), even for **non-normal data**.

22. What Is the Concept of Expected Value in a Probability Distribution?

Expected Value ($E[X]$):

The **expected value** of a random variable is the **long-run average** outcome if the experiment is repeated many times.

Formula:

- For **discrete** random variable:

$$E[X] = \sum (x_i \cdot P(x_i))$$

For **continuous** random variable:

$$E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

Example:

For a fair 6-sided die:

$$E[X] = 1(1/6) + 2(1/6) + \dots + 6(1/6) = 3.5$$

It doesn't mean you'll roll a 3.5, but on **average**, over many rolls, you'll get that value.

Use Cases:

- Decision making
 - Insurance risk calculations
 - Expected profit/loss
-

23. How Does a Probability Distribution Relate to the Expected Outcome of a Random Variable?

Relationship:

A **probability distribution** assigns **probabilities to each possible outcome** of a random variable.

The **expected value** is the **weighted average** of all outcomes based on that distribution.

How They Connect:

- The **distribution** tells us *how likely* each value of the random variable is.
 - The **expected value** summarizes that information into a **single number** representing the **average outcome over time**.
-

Example:

In a game:

- Win ₹100 with 0.3 probability
- Lose ₹50 with 0.7 probability

$$E[X] = (100 \cdot 0.3) + (-50 \cdot 0.7) = 30 - 35 = -₹5$$

So, on average, you lose ₹5 per game.

the End

I done this assignment on google doc

○