

## Statistics Advance-1

Page No.	
Date	

Q.1. Define the Z-statistics and Explain its relationship to the Standard Normal distribution. How is the Z-statistic used in hypothesis testing?

⇒ Z statistic and standard normal distribution

The Z-statistic, or Z-score, measures how many standard deviations a data point (sample mean) is from the population mean. It is a way of standardizing data allowing for comparison between different distributions by converting them to a common scale.

The formula for Z-statistics is:

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

- $x$  is the observed value or sample mean
- $\mu$  is the population mean
- $\sigma$  is the population standard deviation.

If you are dealing with a sample instead of the population, and the population standard deviation is unknown, you would use the sample standard deviation and adjust the formula accordingly when working with the mean of a sample ( $\bar{x}$ ). For a sample population, the formula is..

$$Z = \frac{\bar{x} - \mu}{\sigma}$$

- $\bar{x}$  is the Sample mean.
- $\mu$  is the Population mean.
- $\sigma$  is the Population Standard deviation.
- $n$  is the Sample Size.

### ~~Relationship to the Standard Normal Distribution~~

The  $z$ -statistic follows a standard normal distribution, also known as the  $Z$ -distribution, which has a mean of 0 and a standard deviation of 1. When you calculate a  $Z$ -score by converting raw data into a value on this standard normal distribution, which allows you to assess probabilities and make comparisons.

- A  $Z$ -score of 0 means the value is exactly at the mean.
- Positive  $Z$ -scores indicate values above the mean.
- Negative  $Z$ -scores indicate values below the mean.

The standard normal distribution is essential in statistics because it provides a reference for understanding how extreme or typical a particular point is within a distribution.

### ~~$Z$ statistic in Hypothesis Testing~~

In hypothesis testing, the  $t$  statistic is used when testing the means of large

samples ( $n > 30$ ) or when the population standard deviation is known.

1. Null Hypothesis ( $H_0$ ): A statement that there is no effect or no difference (e.g.  $\mu = \mu_0$ ).
2. Alternative Hypothesis ( $H_A$ ): A statement that there is an effect or a difference (e.g.  $\mu \neq \mu_0$ ).
3. Calculate the Z-statistic, using the formula above to determine how far the sample mean is from the population mean under the null hypothesis.
4. Compare with critical values: the Z-statistic is compared against critical values from the standard normal distribution. Common critical values for a two-tailed test:
  - i.  $|z| > 1.96$  for a 95% confidence level (5% significance level).

→ If  $|z| > 1.96$  you reject the null hypothesis.  
 → If  $|z| \leq 1.96$ , you fail to reject the null hypothesis.

5. P-value: Alternatively, you can use the Z-statistic to calculate the P-value, which tells you the probability of obtaining a test statistic as extreme as the one observed, assuming the null hypothesis is true. If the P-value is less than the significance level (commonly 0.05), you reject the null hypothesis.

### Example

Suppose you want to test out whether the mean height of a population is 170 cm. You collect a sample of 100 people and calculate a sample mean of 172 cm, with a population standard deviation of 10 cm.

Using the Z-statistic formula:

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} = \frac{172 - 170}{10 / \sqrt{100}} = 2$$

with a z-value of 2 and using a significance level of 0.05 (critical z-value is 1.96 for a two-tailed test) you could reject the null hypothesis. Because  $|Z| = 2 > 1.96$ , suggesting that the mean height is significantly different from 170 cm.

Q. 2. What is a p-value and how is it used in hypothesis testing? What does it mean if the p-value is very small (e.g., 0.01)?

⇒ A p-value is the probability of obtaining test results at least as extreme as the observed result under the assumption that the null hypothesis is true. In other words, it quantifies how likely it is to observe the data purely by

random chance if the null hypothesis to hold.

The p-value helps us decide whether the observed data provides enough evidence to reject the null hypothesis in favor of the alternative hypothesis.

### P-value in Hypothesis Testing

In hypothesis testing, the p-value is used to determine whether the observed data is consistent with the null hypothesis.

1. Null Hypothesis ( $H_0$ ): Assumes that there is no effect or no difference. (e.g. the population mean is equal to some hypothesized value).

2. Alternative Hypothesis ( $H_1$ ): Assumes that there is an effect or difference. (e.g. the population mean is different from the hypothesized value).

3. Calculate the Test Statistic: A test statistic (or t-statistic) is calculated based on the sample data.

4. Compute the P-value: The P-value is derived from the test statistic and shows the likelihood of observing the sample data (or more extreme data) if the null hypothesis is true.

5. Compare with Significance Level ( $\alpha$ ): The p-value is compared with predetermined significance level (commonly  $\alpha = 0.05$ ).

→ if  $P\text{-value} \leq \alpha$ : Reject the null hypothesis. This suggests that the observed effect is unlikely under the null hypothesis, providing evidence in favor of the alternative hypothesis.

→ if  $P\text{-value} > \alpha$ : Fail to reject the null hypothesis. This means the data does not provide strongly enough evidence to support the alternative hypothesis.

\* Making out a Small P-value (eg. 0.01):  
if the P-value is very small such as 0.01, it indicates that the observed data is highly unlikely under the null hypothesis. If the null hypothesis is rejected, it means there would be only a 1% chance of obtaining the observed results by random chance. In such cases we conclude that there is strong evidence for the alternative hypothesis.

Interpretation of the final P-value:

P-value = 0.01: There is only 1% probability that the observed effect could occur under the null hypothesis.

Practical interpretation: A small P-value suggests that the observed effect

or difference, is statistically significant. The smaller the p-value, the stronger the evidence against the null hypothesis.

### Example

Suppose you are testing whether a new drug is more effective than the current standard drug. After conducting an experiment, you compute a p-value of 0.01. This means there is only a 1% chance that the observed improvement in patients' outcomes could happen by random chance if the new drug is no better than the standard drug.

→ Since the p-value is smaller than the commonly used significance level of 0.05, you would reject the null hypothesis and conclude that the new drug is significantly more effective.

Q. 3: Compare and contrast the binomial and Bernoulli distributions.

→ The binomial and Bernoulli distributions are closely related, but they are used in different contexts and have some key differences.

#### 1. Definition

\* Bernoulli distribution

→ A Bernoulli distribution models a single trial of a random experiment

With two possible outcomes: "Success" (with probability  $p$ ) and "Failure" (with probability  $1-p$ ).

- it is the simplest discrete probability distribution.
- Example: Flipping a coin once (success = heads, failure = tails).

### ★ Binomial distribution

- A binomial distribution models the number of successes in a fixed number of independent Bernoulli trials.
- It generalizes the Bernoulli distribution by considering multiple trials.
- Example: Flipping a coin 10 times and counting the number of heads.

### 2. Number of Trials

- Bernoulli: A single trial (i.e. one flip of a coin).
- Binomial: involves a fixed number of trials (i.e. flipping the coin 10 times).

### 3. Parameters

Bernoulli: Has only one parameter:

- $p$ : The probability of success.
- The probability of failure is  $1-p$ .
- Binomial: Has two parameters.
- $n$ : the number of independent Bernoulli trials.
- $p$ : the probability of success in each trial.

Page No.	
Date	

## 4. Probability mass function (PMF)

Bernoulli PMF:

$$P(X=0) = p^0(1-p)^1 \quad (62012)$$

where

- $X$  is Bernoulli random variable.
- $x=1$  indicates success (with probability  $p$ ),
- $x=0$  indicates failure (with probability  $1-p$ )

Binomial PMF:

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}, \text{ if } k \in \{0, 1, 2, \dots, n\}$$

where

- $X$  is a binomial random variable representing the number of successes,
- $k$  is the number of successes in trials.
- $\binom{n}{k}$  is the binomial coefficient which counts how many ways you can choose  $k$  successes from  $n$  trials.

## 5. Mean and Variance

Bernoulli Distribution

→ Mean  $\mu = p$

→ Variance  $\sigma^2 = p(1-p)$

Binomial distribution

→ Mean  $\mu = np$

→ Variance  $\sigma^2 = np(1-p)$

## 6. Relationship between the two

→ The Bernoulli distribution is a special case of the binomial distribution where  $n=1$ , in other words, a binomial distribution with  $n=1$  trials is exactly a Bernoulli distribution.

→ Bernoulli: 1 trial with probability  $p$ .  
Binomial:  $n$  independent Bernoulli trials, each with probability  $p$ .

### 7. Example

Bernoulli Example: If a coin has a 50% chance of ~~chance of~~ landing heads ( $P=0.5$ ), a Bernoulli trial is just one flip of the coin. The outcome will either heads (success,  $X=1$ ) or tails (failure,  $X=0$ ).

Binomial Example: If you flip the same coin  $n$  times, the number of heads (successes) will follow a binomial distribution with  $n = 10$  and  $P = 0.5$ . The random variable  $X$  can take values from 0 to 10 (0 heads) to 10 (all heads).

### 8. Use Cases

Bernoulli distribution: Used when modeling the outcome of a single binary event, such as:

- Flipping a coin once
- Whether a student passes or fails an exam.
- Whether a customer buys a product

Binomial distribution: Used when modeling the number of successes in a fixed number of independent binary trials, such as:

- The number of heads in 10 coin flips.

→ The number of students passing can exam out of 50 students.

→ The number of customers making purchases out of 100 who visited a store.

Q.4 Under what conditions is the binomial distribution used, and how does it relate to the Bernoulli distribution?

⇒ The binomial distribution is used under specific conditions when you are dealing with a series of independent trials or experiments, each of which has two possible outcomes: success or failure. It arises in situations where you are interested in the number of successes in a fixed number of trials.

Conditions for using the Binomial distribution:

1. Fixed number of trials ( $n$ ): The experiment or process must consist of a fixed number of trials.

2. Binary outcomes: Each trial must have only two possible outcomes, usually termed "success" or "failure".

3. Constant Probability of Success ( $p$ ): The probability of success on each trial must be the same for every trial.

4. Independence: Each trial must be independent, meaning the outcome of one trial does not affect the outcome of any other trial.

Formula for the Binomial Distribution:

The probability of getting exactly  $k$  successes in  $n$  independent trials is given by:

for the binomial Probability mass function (PMF)

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

where

- $n$  is the number of trials,
- $k$  is the number of successes,
- $p$  is the Probability of success on a single trial and
- $\binom{n}{k}$  is the binomial coefficient, which represents the number of ways to choose  $k$  successes from  $n$  trials.

\* Relationship to the Bernoulli distribution:

⇒ The Bernoulli distribution is a special case of the binomial distribution. It represents the distribution of a single trial (i.e.  $n=1$ ) with two outcomes: success (with probability  $p$ ) or failure (with probability  $1-p$ ).

→ Bernoulli distribution: Describe a single trial with two possible outcomes.

→ Binomial distribution: Describe the number of successes in a fixed number of independent Bernoulli trials.

In other words, the binomial distribution can be thought of as the sum of several independent Bernoulli trials.

QUESTION	
ANSWER	

## Example

Bernoulli trial: Flip a coin once (success = heads, failure = tails).

Binomial Experiment: flip the coin 10 times and count how many heads appear.

→ in this case, each coin flip is a Bernoulli trial and the binomial distribution would describe the probability of getting a specific number of heads (success) in the flips (trials).

Q.5. What are the key properties of the Poisson distribution and when is it appropriate to use this distribution?

⇒ The Poisson distribution is a probability used to model the number of events that occur in a fixed set of intervals of time, space or another domain, under specific conditions. It is particularly useful for modeling rare or infrequent events.

Key properties of the Poisson distribution:

1. Discrete distribution: The Poisson distribution applies to discrete events (i.e. events that can only take on non-negative integer values: 0, 1, 2, ...).

2. Rate of occurrence ( $\lambda$ ): the parameter  $\lambda$  represents both the average number of events in the given interval and the rate at which the events occur; it must be positive.

3. Independence: The occurrence of one event does not affect the occurrence of another event. Each event happens independently.
4. Events occur Randomly: Events are distributed randomly over time.
5. Non-overlapping intervals: The number of events occurring in one interval is independent of the number of events occurring in any other non-overlapping interval.
6. No upper limit: Theoretically, the number of events can range from zero to infinity, though the probability of very high counts decreases as the count increases.

7. Mean and Variance: The mean (Expected value) and variance of a Poisson distribution are both equal to  $\lambda$ .

- Mean:  $E(X) = \lambda$
- Variance:  $Var(X) = \lambda$

### Poisson Distribution Formula

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

where

- $\rightarrow X$  is the random variable representing the number of events
- $\rightarrow k$  is the number of events (0, 1, 2, ...)
- $\rightarrow \lambda$  is the average rate (mean number of events) in the given interval.

→ e is the base of the natural logarithm (approximately 2.718)

\* when to use the Poisson distribution

1. modeling the number of events in a fixed interval of time, space, or volume. P.e.g., number of phone call received in an hour, or number of customer arrivals at a store, or number of earth quakes in year.
2. The events are rare and occurs infrequently relative to the length of the interval.
3. the events occur independently of each other.
4. the average rate of occurrence is constant over time or space. (i.e.  $\lambda$  remains constant).
5. The probability of more than one event occurring in a very short interval is negligible.

### Use Cases

Call center, Accidents, Arrivals, Natural Phenomena.

Q. Define the terms "Probability distribution" and "Probability density function" (PDF). How does a PDF differ a Probability mass Function (PMF)?

⇒ probability distribution

A probability distribution is a mathematical description of the likelihood of different outcomes in an experiment or random

process. It assigns a probability to each possible outcome of a random variable, depending on whether the random variable is discrete or continuous.

- For discrete random variables, the distribution is described by a probability mass function.
- For continuous random variables, the distribution is described by a probability density function (PDF).

A probability distribution provides a complete description of a random variable's behavior, including the possible values it can take and the likelihood (or density) of these values.

\* Probability density function (PDF)  
A probability density function (PDF) describes the likelihood of a continuous random variable taking a particular value. Unlike a discrete probability distribution where probabilities are assigned to specific outcomes, the PDF gives the density or probability at any given value of the random variable.

- A PDF,  $f_{X(x)}$ , for a continuous random variable,  $X$ , satisfies the following properties:
1.  $f_{X(x)} \geq 0$  for all  $x$ .
  2. The total area under the curve is 1.

for PDF (which is 1 i.e.  $\int_{-\infty}^{\infty} f(x) dx = 1$ ).

3. The probability that the variable  $X$  falls within a particular range  $[a, b]$  is given by the area under the curve between  $a$  and  $b$ :

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

In Continuous distribution, the probability of the random variable taking on an exact value (e.g.  $P(X=2)$ ) is zero, but probabilities for intervals (e.g.  $P(a \leq X \leq b)$ ) can be calculated.

### 4. Probability Mass Function (PMF)

A Probability mass function (PMF) is used for discrete random variables and gives the probability that the random variable takes on a specific value.

→ A PMF,  $P(X=x)$ , satisfies:

- $P(X=x) \geq 0$  for all  $x$ .

2. The sum of probabilities across all possible values of  $X$  equals 1:

$$\sum_x P(X=x) = 1$$

In a Discrete Distribution, the PMF assigns a specific probability to each possible outcome.

Differences between PDF and PMF

\* Type of random variable:

PMF is used for discrete random

→ variables where outcomes take distinct, countable values (e.g. 0, 1, 2, ...)

→ PDF is used for continuous random variables, where outcomes can take any value within an interval (e.g. any real number).

\* Interpretation of Values:

→ for a PMF, the function directly gives the probability of the random variable taking on a specific value (e.g.  $P(X = 0.3)$ ).

→ for a PDF, the function gives the density of the random variable at a point, not the actual probability. Probabilities are obtained by integrating the PDF over an interval. (e.g.,  $P(a \leq X \leq b)$ ).

\* Simulation vs. Integration

→ for PMFs, probabilities are summed over possible values.

→ For PDFs, probabilities are calculated by integrating over an interval of values. Example:

PMF (Discrete case): Consider a fair six-sided die. The PMF is:

$$P(X=k) = \frac{1}{6} \text{ for } k = 1, 2, 3, 4, 5, 6$$

Each outcome has an equal probability of  $\frac{1}{6}$ .

Q.2. PDF (Continuous Case): For a random variable that is normally distributed with mean  $\mu$  and Variance  $\sigma^2$ , the PDF is:

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

Q.7. Explain the Central Limit Theorem (CLT) with example.

$\Rightarrow$  The Central limit theorem (CLT) is one of the most important results in statistics. It states that, regardless of the shape of the original population distribution, the distribution of the sample mean (or the sum of a sufficiently large number of independent and identically distributed random variables will tend to follow a normal distribution (a bell curve) provided the sample size is large enough.

- \* Key Points of the Central Limit Theorem
  1. Independence: The random variable being considered or summed must be independent.
  2. Identically Distributed: The random variables should come from the same probability distribution.
  3. Sample size: The sample size must be sufficiently large ( $n \geq 30$ ) is considered large enough.
  4. Mean and Variance: The mean of the sample mean distribution is equal to the mean of the original distribution and the variance of the sample mean distribution is equal to the

Variance at the original distribution divided by  $n$  (the sample size).

#### \* Mathematical Statement

Let  $X_1, X_2, \dots, X_n$  be i.i.d random variables from a population with mean  $\mu$  and variance  $\sigma^2$ . The CLT states that the distribution of the sample mean  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  approaches a normal distribution with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$  as  $n$  becomes large. In other words,

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \text{ as } n \rightarrow \infty$$

#### \* Example of Central Limit Theorem

Suppose you are rolling a die. The outcomes of each roll follows a uniform distribution (i.e. Each face has an equal probability of  $\frac{1}{6}$ ). The mean  $\mu$  of this distribution is 3.5 and the variance  $\sigma^2$  can be calculated for this uniform distribution.

Now, let's apply the CLT.

1. Single Roll: The distribution of outcome at a single roll is not normal - it's uniform.
2. Average of multiple rolls: If you roll the die 30 times and compute the average outcome over those 30 rolls,

The distribution of the average will be approximately normal even though the individual roll outcomes are not normally distributed.

→ As the number of rolls increases, the distribution of the sample mean (average outcome) becomes closer and closer to a normal distribution with mean  $\mu = 3.5$  and standard deviation  $\sigma$ , where  $\sigma$  is the standard deviation of the uniform distribution.

\* Visualization of the central limit theorem:

1. If you roll a die 1 time, the distribution is flat (uniform).
2. If you roll a die 10 times and plot the average outcome from many experiments of 10 rolls, the distribution starts resembling a bell curve.
3. If you roll a die 50 times and plot the average outcome from many experiments of 50 rolls, the distribution of those averages is even closer to a normal distribution.

\* importance of the central limit theorem

→ Real World Application: The ability to use normal distribution techniques even when the underlying data is not normally distributed, as long as one is carefully working with sample means or sums.

→ Statistical inference: The CLT Justifies the use of normal distribution-based methods (Confidence Intervals and hypothesis testing) in many real-world situations; even when the population distribution is unknown or non-normal.

Q. 8. Compare Z-scores and t-scores. When should you use a Z-score and when should a t-score be applied instead?

⇒ Z-score and t-scores (or t-values) are both used in statistics to measure how far a data point or a sample statistic (like the sample mean) is from the population mean, measured in terms of standard deviations or standard errors. However, they are used in different contexts depending on the available data and certain assumptions about the population.

### Z-score

A Z-score measures how many standard deviations an individual data point or sample statistic is away from the population mean if it assumes that the population standard deviation  $\sigma$  is known and that the population is normally distributed.

→ Formula for Z-score:

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

where

- $x$  is the mean or sample mean,
- $\mu$  is the population mean.
- $\sigma$  is the population standard deviation.

### \* When to use $Z$ -score

1. Large sample size: when the sample size  $n$  is large (typically  $n \geq 30$ ).
2. Known Population Standard Deviation ( $\sigma$ ): when you know the population standard deviation.
3. Normal distribution: Z-scores are often used when the population is normally distributed; but due to the Central Limit theorem they can also be applied when  $n \geq 30$  even if the population distribution is not normal.

### \* Z-Score Rules

- Finding probabilities using the standard normal distribution (Z-table) for standardized test scores.
- Calculating confidence intervals or conducting hypothesis tests when the population standard deviation is known.

### \* T-Score:

A t-score (or t-value) is used when the sample size is small and the population standard deviation  $\sigma$  is unknown; instead of using the population standard deviation  $\sigma$ , we use the sample standard deviation  $s$ ; which adds variability and uncertainty to the estimate. The distribution is

Similar to the normal distribution but has heavier tails, meaning there is more probability in the extremes due to the uncertainty in estimating the standard deviation

t-score formula:

$$t = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$$

where

- $\bar{X}$  is the sample mean
- $\mu$  is the population mean
- $s$  is the sample standard deviation
- $n$  is the sample size

\* when to use a t-score

1. small sample size: when the sample size,  $n$ , is small (typically  $n < 30$ )
2. unknown population standard deviation: when the population standard deviation  $\sigma$  is unknown and you are using the sample standard deviation  $s$  as an estimate.

3. nearly normal distribution: the t-distribution assumes the sample comes from a normally or nearly normally distributed population, especially for small sample sizes. As  $n$  increases, the t-distribution approaches the normal distribution.

## \* Examples

- Calculating confidence intervals or conducting hypothesis tests for small sizes where  $\sigma$  is unknown.
- Student's t-test for comparing means of small samples sizes where  $\sigma$  is unknown.

## \* Example of when to use each

### 1. Z-score Example:

- Suppose you want to calculate a confidence interval for the mean of a large sample of 100 test scores where the population standard deviation  $\sigma$  is known. You would use the Z-score to construct the interval.

### 2. T-Score Example

- Suppose you have a small sample of 15 students test scores, and you don't know the population standard deviation. You would use the t-score to calculate the confidence interval, accounting for the added uncertainty in estimating the population standard deviation.