

Statistical Inference Foundations

Vaibhav P. Vasani

Assistant Professor

Department of Computer Engineering

K. J. Somaiya College of Engineering

Somaiya Vidyavihar University

Statistical Inference

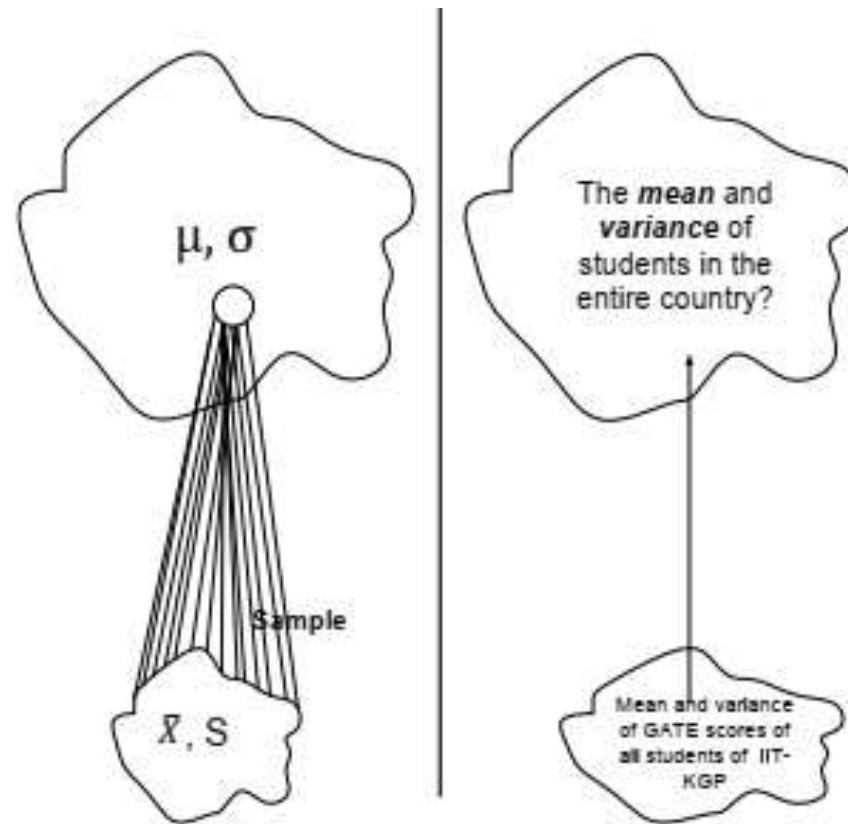
- Using data analysis and statistics to make conclusions about a population is called statistical inference.

The main types of statistical inference are:

- Estimation
- Hypothesis testing



The primary objective of statistical analysis is to use data from a sample to make inferences about the population from which the sample was drawn.





SOMAIYA
VIDYAVIHAR UNIVERSITY

K J Somaiya College of Engineering





SOMAIYA
VIDYAVIHAR UNIVERSITY

K J Somaiya College of Engineering



Estimation

- Statistics from a sample are used to estimate population parameters.
- The most likely value is called a **point estimate**.
- There is **always** uncertainty when estimating.
- The uncertainty is often expressed as **confidence intervals** defined by a likely lowest and highest value for the parameter.



- An example could be a confidence interval for the number of bicycles a Dutch person owns:
 - "The average number of bikes a Dutch person owns is between 3.5 and 6."



Hypothesis Testing

- **Hypothesis testing** is a method to check if a claim about a population is true. More precisely, it checks how likely it is that a hypothesis is true is based on the sample data.
- There are different types of hypothesis testing.
- The steps of the test depends on:
- Type of data (categorical or numerical)
- If you are looking a at:
 - A single group
 - Comparing one group to another
 - Comparing the same group before and after a change

- Some examples of claims or questions that can be checked with hypothesis testing:
 - 90% of Australians are left handed
 - Is the average weight of dogs more than 40kg?
 - Do doctors make more money than lawyers?



Note

- Probability Distributions
- Statistical inference methods rely on probability calculation and probability distributions.



The 30 homes in this dataset have been selected randomly from the population of all single-family homes for sale in this housing market. We can simply list small datasets such as this. The values of Price in this case are:

155.5	195.0	197.0	207.0	214.9	230.0	239.5	242.0	252.5	255.0
259.9	259.9	269.9	270.0	274.9	283.0	285.0	285.0	299.0	299.9
319.0	319.9	324.5	330.0	336.0	339.0	340.0	355.0	359.9	359.9



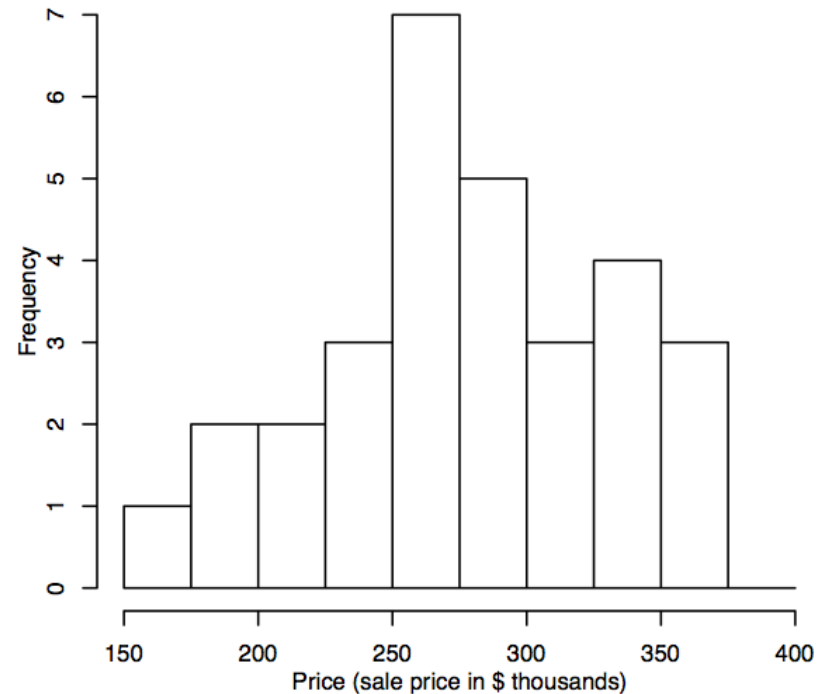
- A particularly effective graph here is a stem-and-leaf plot, which places the numbers along the vertical axis of the plot, with numbers that are close together in magnitude next to one another on the plot. For example, a stem-and-leaf plot for the 30 sample prices looks like the following:

- 1 | 6
- 2 | 0011344
- 2 | 5666777899
- 3 | 002223444
- 3 | 666

- In this plot, the decimal point is two digits to the right of the stem. So, the “1” in the stem and the “6” in the leaf represents 160, or, because of rounding, any number between 155 and 164.9. In particular, it represents the lowest price in the dataset of 155.5 (thousand dollars). The next part of the graph shows two prices between 195 and 204.9, two prices between 205 and 214.9, one price between 225 and 234.9, two prices between 235 and 244.9, and so on. A stem-and-leaf plot can easily be constructed by hand for small datasets such as this, or it can be constructed automatically using statistical software.



- A few modifications to a stem-and-leaf plot produce a histogram—the value axis is now horizontal rather than vertical, and the counts of observations within adjoining data intervals (called “bins”) are displayed in bars (with the counts, or frequency, shown on the vertical axis) rather than by displaying individual values with digits. The following shows a histogram for the home prices data generated by statistical software.



•



SOMAIYA
VIDYAVIHAR UNIVERSITY

K J Somaiya College of Engineering



Basic Approaches

Approach 1: Hypothesis testing

- ▶ We conduct **test on hypothesis**.
 - ▶ We hypothesize that one (or more) parameter(s) has (have) some specific value(s) or relationship.
- ▶ Make our decision about the parameter(s) based on one (or more) sample statistic(s)
- ▶ Accuracy of the decision is expressed as the probability that the **decision is incorrect**.

Approach 2: Confidence interval measurement

- ▶ We estimate one (or more) parameter(s) using sample statistics.
 - ▶ This estimation usually done in the form of an interval.
- ▶ Accuracy of the decision is expressed as the **level of confidence** we have in the interval.

Hypothesis Testing

- **What is Hypothesis?**
- “A hypothesis is an educated prediction that can be tested” ([study.com](#)).
- “A hypothesis is a proposed explanation for a phenomenon” ([Wikipedia](#)).
- “A hypothesis is used to define the relationship between two variables” ([Oxford dictionary](#)).
- “A supposition or proposed explanation made on the basis of limited evidence as a starting point for further investigation” ([Walpole](#)).

Example 6.1: Avogadro's Hypothesis(1811)

“The volume of a gas is directly proportional to the number of molecules of the gas.”



Statistical Hypothesis

- ▶ If the hypothesis is stated in terms of population parameters (such as mean and variance), the hypothesis is called **statistical hypothesis**.
- ▶ Data from a sample (which may be an experiment) are used to test the validity of the hypothesis.
- ▶ A procedure that enables us to agree (or disagree) with the statistical hypothesis is called a **test of the hypothesis**.

Example 6.2:

1. To determine whether the wages of men and women are equal.
2. A product in the market is of standard quality.
3. Whether a particular medicine is effective to cure a disease.

The Hypotheses

- ▶ The main purpose of statistical hypothesis testing is to choose between two competing hypotheses.

Example 6.3:

One hypothesis might claim that wages of men and women are equal, while the **alternative** might claim that men make more than women.

- ▶ Hypothesis testing start by making a set of two statements about the parameter(s) in question.
- ▶ The hypothesis actually to be tested is usually given the symbol H_0 and is commonly referred as the **null hypothesis**.
- ▶ The other hypothesis, which is assumed to be true when null hypothesis is false, is referred as the **alternate hypothesis** and is often symbolized by H_1
- ▶ The two hypotheses are **exclusive** and **exhaustive**.

Hypothesis Testing Procedures

The following **five steps** are followed when testing hypothesis

1. Specify H_0 and H_1 , the null and alternate hypothesis, and an **acceptable level of α** .
2. Determine an appropriate sample-based test statistics and the **rejection region** for the specified H_0 .
3. Collect the sample data and calculate the test statistics.
4. Make a decision to either reject or fail to reject H_0 .
5. Interpret the result in common language suitable for practitioners.



Hypothesis Testing Procedure

- ▶ In summary, we have to choose between H_0 and H_1
- ▶ The standard procedure is to assume H_0 is true.

(Just we presume innocent until proven guilty)
- ▶ Using statistical test, we try to determine whether there is sufficient evidence to declare H_0 false.
- ▶ We reject H_0 only when the **chance is small** that H_0 is true.
- ▶ The procedure is based on probability theory, that is, there is a chance that we can **make errors**.

Errors in Hypothesis Testing

- In hypothesis testing, there are two types of errors.

Type I error: A type I error occurs when we incorrectly reject H_0 (i.e., we reject the null hypothesis, when H_0 is true).

Type II error: A type II error occurs when we incorrectly fail to reject H_0 (i.e., we accept H_0 when it is not true).

Decision	Observation	
	H_0 is true	H_0 is false
H_0 is accepted	Decision is correct	Type II error
H_0 is rejected	Type I error	Decision is correct

Probabilities of Making Errors

Type I error calculation

α denotes the probability of making a Type I error

$$\alpha = P(\text{Rejecting } H_0 | H_0 \text{ is true})$$

Type II error calculation

β denotes the probability of making a Type II error

$$\beta = P(\text{Accepting } H_0 | H_0 \text{ is false})$$

Note:

- ▶ α and β are not independent of each other as one increases, the other decreases
- ▶ When the sample size increases, both to decrease since sampling error is reduced.
- ▶ In general, we focus on Type I error, but Type II error is also important, particularly when sample size is small.

Calculating α

Assuming that we have the results of random sample. Hence, we use the characteristics of sampling distribution to calculate the probabilities of making either Type I or Type II error.

Example 6.6:

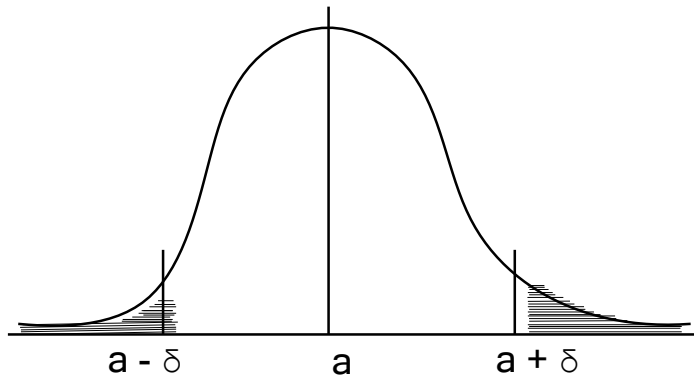
Suppose, two hypotheses in a statistical testing are:

$$H_0: \mu = a$$

$$H_1: \mu \neq a$$

Also, assume that for a given sample, population obeys normal distribution. A threshold limit say $a \pm \delta$ is used to say that **they are significantly different from a .**

Calculating α



Here, shaded region implies the probability that $\bar{X} < a - \delta$ or $\bar{X} > a + \delta$

Thus the null hypothesis is to be rejected if the mean value is less than $a - \delta$ or greater than $a + \delta$.

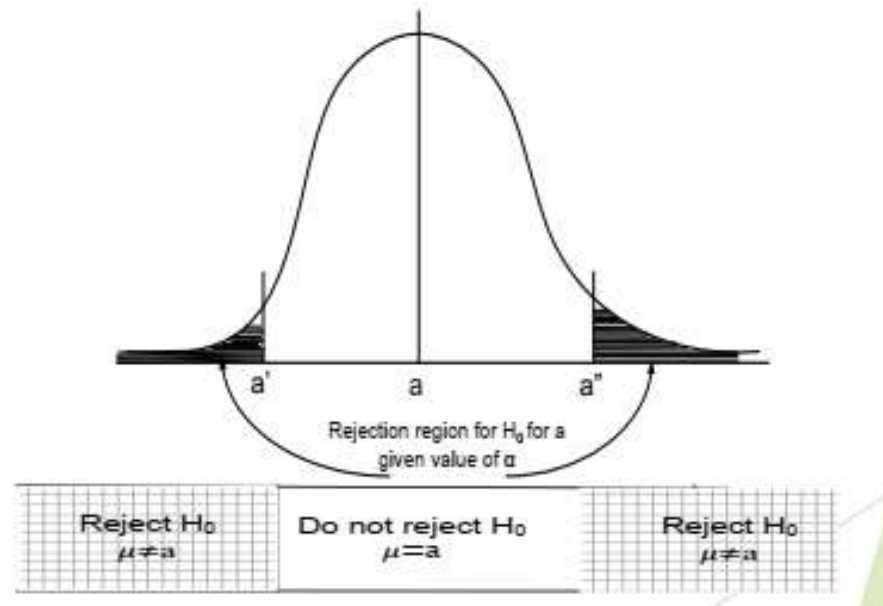
If \bar{X} denotes the sample mean, then the Type I error is

$$\alpha = P(\bar{X} < a - \delta \text{ or } \bar{X} > a + \delta, \text{ when } \mu = a, \text{ i.e., } H_0 \text{ is true})$$

The Rejection Region

The rejection region comprises of value of the test statistics for which

1. The probability when the null hypothesis is true is less than or equal to the specified α .
2. Probability when H_1 is true are greater than they are under H_0 .



Two-Tailed Test

For two-tailed hypothesis test, hypotheses take the form

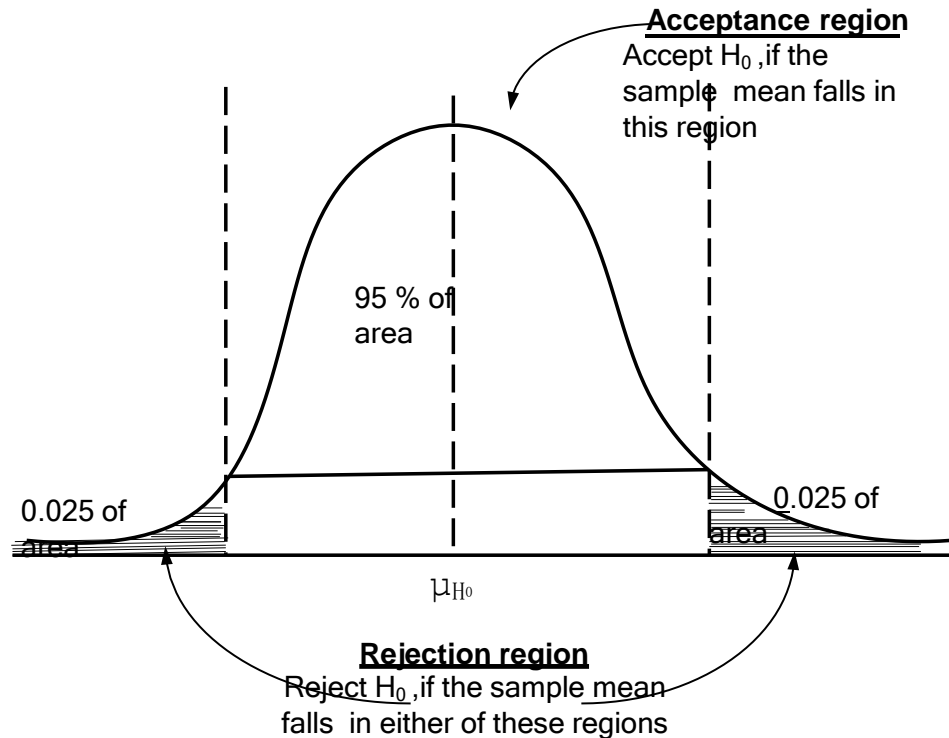
$$H_0: \mu = \mu_{H_0}$$

$$H_1: \mu \neq \mu_{H_0}$$

In other words, to reject a null hypothesis, sample mean $\mu > \mu_{H_0}$ or $\mu < \mu_{H_0}$ under a given α .

Thus, in a two-tailed test, there are two rejection regions (also known as critical region), one on each tail of the sampling distribution curve.

Two-Tailed Test



Acceptance and rejection regions in case of a two-tailed test with 5% significance level.

One-Tailed Test

A one-tailed test would be used when we are to test, say, whether the population mean is either lower or higher than the hypothesis test value.

Symbolically

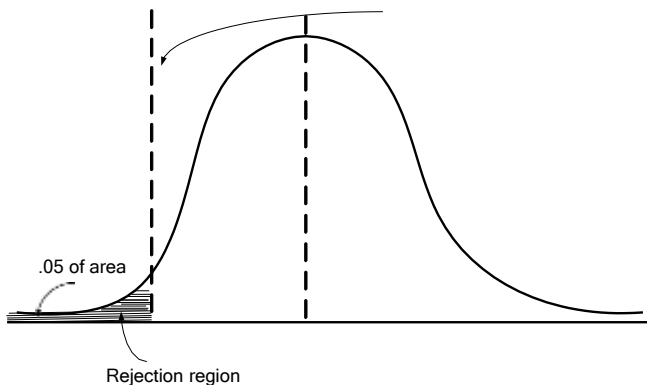
$$H_0: \mu = \mu_{H_0}$$

,

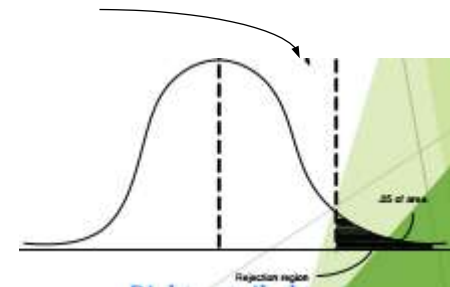
$$H_1: \mu < \mu_{H_0}$$

$$[or \mu > \mu_{H_0}]$$

Wherein there is one rejection region only on the left-tail (or right tail)



Left - tailed



Right - tailed
test

Example 6.7: Calculating α

Consider the two hypotheses are

The null hypothesis is

$$H_0: \mu = 8$$

The alternative hypothesis is

$$H_1: \mu \neq 8$$

Assume that given a sample of size 16 and standard deviation is 0.2 and sample follows normal distribution.

Example 6.7: Calculating α

We can decide the rejection region as follows.

Suppose, the null hypothesis is to be rejected if the mean value is less than 7.9 or greater than 8.1. If \bar{X} is the sample mean, then the probability of Type I error is

$$\alpha = P(\bar{X} < 7.9 \text{ or } \bar{X} > 8.1, \text{ when } \mu = 8)$$

Given σ , the standard deviation of the sample is 0.2 and that the distribution follows **normal distribution**.

Thus,

$$P(\bar{X} < 7.9) = P\left[Z = \frac{7.9 - 8}{0.2/\sqrt{16}}\right] = P[Z < -2.0] = 0.0228$$

and

$$P(\bar{X} > 8.1) = P\left[Z = \frac{8.1 - 8}{0.2/\sqrt{16}}\right] = P[Z > 2.0] = 0.0228$$

Hence, $\alpha = 0.0228 + 0.0228 = 0.0456$

Example 6.8: Calculating α and β

There are two identically appearing boxes of chocolates. Box A contains 60 red and 40 black chocolates whereas box B contains 40 red and 60 black chocolates. There is no label on the either box. One box is placed on the table. We are to test the hypothesis that “Box B is on the table”.

To test the hypothesis an experiment is planned, which is as follows:

- ▶ Draw at random five chocolates from the box.
- ▶ We replace each chocolates before selecting a new one.
- ▶ The number of red chocolates in an experiment is considered as the **sample statistics**.

Note: Since each draw is independent to each other, we can assume the sample distribution follows binomial probability distribution.

Example 6.8: Calculating α

Let us express the population parameter as p = the number of red chocolates in Box B.

The hypotheses of the problem can be stated as:

$$H_0: p = 0.4 \quad // \text{ Box B is on the table}$$

$$H_1: p = 0.6 \quad // \text{ Box A is on the table}$$

Calculating α :

In this example, the null hypothesis (H_0) specifies that the probability of drawing a red chocolate is 0.4. This means that, lower proportion of red chocolates in observations (*i. e.*, *sample*) favors the null hypothesis. In other words, **drawing all red chocolates** provides **sufficient evidence to reject the null hypothesis**. Then, the probability of making a *Type I* error is the probability of getting five red chocolates in a sample of five from Box B. That is,

$$\alpha = P(X = 5 \text{ when } p = 0.4)$$

Using the binomial distribution

$$\begin{aligned} &= \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \text{ where } n = 5, x = 5 \\ &= (0.4)^5 = 0.01024 \end{aligned}$$

Thus, the probability of rejecting a true null hypothesis is ≈ 0.01 . That is, there is approximately 1 in 100 chance that the box B will be mislabeled as box A.

Example 6.8: Calculating β

The *Type II* error occurs if we fail to reject the null hypothesis when it is not true. For the current illustration, such a situation occurs, if Box A is on the table but we did not get the five red chocolates required to reject the hypothesis that Box B is on the table.

The probability of *Type II* error is then the probability of getting four or fewer red chocolates in a sample of five from Box A.

That is,

$$\beta = P(X \leq 4) \quad \text{when } p = 0.6)$$

Using the probability rule:

$$P(X \leq 4) + P(X = 5) = 1$$

$$\text{That is, } P(X \leq 4) = 1 - P(X = 5)$$

$$\text{Now, } P(X = 5) = (0.6)^5$$

$$\begin{aligned} \text{Hence, } \beta &= 1 - (0.6)^5 \\ &= 1 - 0.07776 = 0.92224 \end{aligned}$$

That is, the probability of making *Type II* error is over 92%. This means that, if Box A is on

Case Study 1: Coffee Sale

A coffee vendor nearby Kharagpur railway station has been having average sales of 500 cups per day. Because of the development of a bus stand nearby, it expects to increase its sales. During the first 12 days, after the inauguration of the bus stand, the daily sales were as under:

550 570 490 615 505 580 570 460 600 580 530 526

On the basis of this sample information, can we conclude that the sales of coffee have increased?

Consider 5% level of confidence.



Hypothesis Testing : 5 Steps

The following **five steps** are followed when testing hypothesis

1. Specify H_0 and H_1 , the null and alternate hypothesis, and an **acceptable level of α** .
2. Determine an appropriate sample-based test statistics and the **rejection region** for the specified H_0 .
3. Collect the sample data and calculate the test statistics.
4. Make a decision to either reject or fail to reject H_0 .
5. Interpret the result in common language suitable for practitioner.



Case Study 1: Step 1

Step 1: Specification of hypothesis and acceptable level of α

Let us consider the hypotheses for the given problem as follows.

$$H_0: \mu = 500 \text{ cups per day}$$

The null hypothesis that sales average 500 cups per day and they have not increased.

$$H_1: \mu > 500$$

The alternative hypothesis is that the sales have increased.

Given the acceptance level of $\alpha = 0.05$ (i. e. , 5% level of significance)

Case Study 1: Step 2

Step 2: Sample-based test statistics and the rejection region for specified H_0

Given the sample as

550 570 490 615 505 580 570 460 600 580 530 526

Since the sample size is small and the population standard deviation is not known, we shall use

t - test assuming normal population. The test statistics t is

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

To find \bar{X} and S , we make the following computations.

$$\bar{X} = \frac{\sum x_i}{n} = \frac{6576}{12} = 548$$

Case Study 1: Step 2

Sample #	X_i	$X_i - \bar{X}$	$(X_i - \bar{X})^2$
1	550	2	4
2	570	22	484
3	490	-58	3364
4	615	67	4489
5	505	-43	1849
6	580	32	1024
7	570	22	484
8	460	-88	7744
9	600	52	2704
10	580	32	1024
11	530	-18	324
12	526	-22	484
$n = 12$	$\sum X_i = 6576$		$\sum (X_i - \bar{X})^2 = 23978$

Case Study 1: Step 2

$$S = \sqrt{\frac{\sum (X_i - \bar{X})^2}{n - 1}} = \sqrt{\frac{23978}{12 - 1}} = 46.68$$

$$\text{Hence, } t = \frac{\bar{X} - \mu}{S/\sqrt{n}} = \frac{48}{46.68/\sqrt{12}} = \frac{48}{13.49} = 3.558$$

Note:

Statistical table for t-distributions gives a t -value given n , the degrees of freedom and α , the level of significance and vice-versa.

Case Study 1: Step 3

Step 3: Collect the sample data and calculate the test statistics

$$\text{Degree of freedom} = n - 1 = 12 - 1 = 11$$

As H_1 is one-tailed, we shall determine the rejection region applying one-tailed in the right tail because H_1 is more than type) at 5% level of significance.



Case Study 1: Step 3

Step 3: Collect the sample data and calculate the test statistics

$$\text{Degree of freedom} = n - 1 = 12 - 1 = 11$$

As H_1 is one-tailed, we shall determine the rejection region applying one-tailed in the right tail because H_1 is more than type) at 5% level of significance.

Using table of t - *distribution* for 11 degrees of freedom and with 5% level of significance,

$$R: t > 1.796$$

Case Study 1: Step 4

Step 4: Make a decision to either reject or fail to reject H_0

The observed value of $t = 3.558$ which is in the rejection region and thus H_0 is rejected at 5% level of significance.

Case Study 1: Step 5

Step 5: Final comment and interpret the result

We can conclude that the sample data indicate that coffee sales have increased.



Case Study 2: Machine Testing

A medicine production company packages medicine in a tube of 8 ml. In maintaining the control of the amount of medicine in tubes, they use a machine. To monitor this control a sample of 16 tubes is taken from the production line at random time interval and their contents are measured precisely. The mean amount of medicine in these 16 tubes will be used to test the hypothesis that the machine is indeed working properly.



Case Study 2: Step 1

Step 1: Specification of hypothesis and acceptable level of α

The hypotheses are given in terms of the population mean of medicine per tube.

The null hypothesis is

$$H_0: \mu = 8$$

The alternative hypothesis is

$$H_1: \mu \neq 8$$

We assume α , the significance level in our hypothesis testing ≈ 0.05 .

(This signifies the probability that the machine needs to be adjusted less than 5%).

Case Study 2: Step 2

Step 2: Sample-based test statistics and the rejection region for specified H_0

Rejection region: Given $\alpha = 0.05$, which gives $|Z| > 1.96$ (obtained from standard normal calculation for $n(Z: 0, 1) = 0.025$ for a rejection region with two-tailed test).

Case Study 2: Step 3

Step 3: Collect the sample data and calculate the test statistics

Sample results: $n = 16$, $\bar{x} = 7.89$, $\sigma = 0.2$

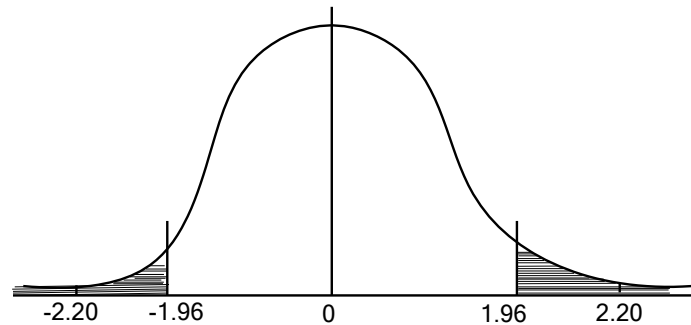
With the sample, the test statistics is

$$Z = \frac{7.89 - 8}{\frac{0.2}{\sqrt{16}}} = -2.20$$

Hence, $|Z| = 2.20$

Case Study 2: Step 4

Step 4: Make a decision to either reject or fail to reject H_0



Since $Z > 1.96$, we reject H_0

Case Study 2: Step 5

Step 5: Final comment and interpret the result

We conclude $\mu \neq 8$ and recommend that the machine be adjusted.



Case Study 2: Alternative Test

Suppose that in our initial setup of hypothesis test, if we choose $\alpha = 0.01$ instead of 0.05, then the test can be summarized as:

1. $H_0: \mu = 8$, $H_1: \mu \neq 8$ $\alpha = 0.01$
2. Reject H_0 if $Z > 2.576$
3. Sample result $n=16$, $\sigma = 0.2$, $\bar{X}=7.89$, $Z = \frac{7.89-8}{\frac{0.2}{\sqrt{16}}} = -2.20$, $|Z| = 2.20$
4. $|Z| < 2.20$, we fail to reject $H_0=8$
5. We do not recommend that the machine be readjusted.



Hypothesis Testing Strategies

- ▶ The hypothesis testing determines the validity of an assumption (technically described as null hypothesis), with a view to choose between two conflicting hypothesis about the value of a **population** parameter.
- ▶ There are two types of tests of hypotheses
 - ✓ Non-parametric tests (also called distribution-free test of hypotheses)
 - ✓ Parametric tests (also called standard test of hypotheses).

Parametric Tests : Applications

- ▶ Usually assume certain properties of the population from which we draw samples.
 - Observation come from a normal population
 - Sample size is small
 - Population parameters like mean, variance, etc. are hold good.
 - Requires measurement equivalent to interval scaled data.

Parametric Tests

Important Parametric Tests

The widely used sampling distribution for parametric tests are

- ▶ Z – test
- ▶ t – test
- ▶ χ^2 – test
- ▶ F – test

Note:

All these tests are based on the assumption of normality (i.e., the source of data is considered to be normally distributed).

Parametric Tests : Z-test

Z – test: This is most frequently test in statistical analysis.

- ▶ It is based on the normal probability distribution.
- ▶ Used for judging the significance of several statistical measures particularly the mean.
- ▶ It is used even when *binomial distribution* or *t – distribution* is applicable with a condition that such a distribution tends to normal distribution when n becomes large.
- ▶ Typically it is used for comparing the mean of a sample to some hypothesized mean for the population in case of large sample, or when **population variance** is known.

Parametric Tests : t-test

***t* – *test*:** It is based on the t-distribution.

- ▶ It is considered an appropriate test for judging the significance of a sample mean or for judging the significance of difference between the means of two samples in case of
 - ▶ small sample(s)
 - ▶ **population variance is not known** (in this case, we use the variance of the sample as an estimate of the population variance)



Parametric Tests : χ^2 -test

χ^2 – *test*: It is based on Chi-squared distribution.

- ▶ It is used for comparing a sample variance to a theoretical population variance.

Parametric Tests : F -test

F – *test*: It is based on F-distribution.

- ▶ It is used to compare the variance of two *independent samples*.
- ▶ This test is also used in the context of analysis of variance (ANOVA) for judging the significance of more than two sample means.

Hypothesis Testing : Assumptions

Case 1: Normal population, population infinite, sample size may be large or small, variance of the population is known.

$$z = \frac{\bar{X} - \mu_{H_0}}{\sigma/\sqrt{n}}$$

Case 2: Population normal, population **finite**, sample size may large or small.....variance is known.

$$z = \frac{\bar{X} - \mu_{H_0}}{\sigma/\sqrt{n}[\sqrt{(N-n)/(N-1)}]}$$

Case 3: Population normal, population infinite, **sample size is small** and variance of the **population is unknown**.

$$t = \frac{\bar{X} - \mu_{H_0}}{s/\sqrt{n}} \\ (n-1)$$

with degree of freedom =

and

$$s = \sqrt{\frac{\sum(X_i - \bar{X})^2}{(n-1)}}$$

Hypothesis Testing

Case 4: Population finite

$$t = \frac{\bar{X} - \mu_{H_0}}{\sigma / \sqrt{n} [\sqrt{(N-n)/(N-1)}]} \text{ with degree of freedom} = (n - 1)$$

Note: If variance of population (σ) is known, replace S by σ . Population normal, population infinite, **sample size is small** and variance of the **population is unknown**.

Hypothesis Testing : Non-Parametric Test

► *Non-Parametric tests*

- ✓ Does not under any assumption
- ✓ Assumes only nominal or ordinal data

Note: Non-parametric tests need entire population (or very large sample size)



SOMAIYA
VIDYAVIHAR UNIVERSITY

K J Somaiya College of Engineering



Reference

- The detail material related to this lecture can be found in

Probability and Statistics for Engineers and Scientists (8th Ed.) by Ronald E. Walpole, Sharon L. Myers, Keying Ye (Pearson), 2013.

Question ?



SOMAIYA
VIDYAVIHAR UNIVERSITY

K J Somaiya College of Engineering

