# Statistics for Data Analysis-Lec 4

**Lecturer:** **Taufique Ahmed**

**E-mail:** **tahmed@cct.ie**

cct | College Dublin
Computing • IT • Business

College Dublin
Computing • IT • Business

# Hypothesis Testing

# Hypothesis Testing

- Hypothesis testing can be used to determine whether a statement about the value of a population parameter should or should not be rejected.

- The null hypothesis, denoted by $H_0$, is a tentative assumption about a population parameter

- The alternative hypothesis, denoted by Ha, is the opposite of what is stated in the null hypothesis

- The hypothesis testing procedure uses data from a sample to test the two competing statements indicated by H0 and Ha.

# Developing Null and Alternative Hypotheses

- It is not always obvious how the null and alternative hypotheses should be formulated

- Care must be taken to structure the hypotheses appropriately so that the test conclusion provides the information the researcher wants

- The context of the situation is very important in determining how the hypotheses should be stated

- In some cases it is easier to identify the alternative hypothesis first. In order cases the null is easier

- Correct hypothesis formulation will take practice

# Developing Null and Alternative Hypotheses

**Alternative Hypothesis as a Research Hypothesis**

- Many applications of hypothesis testing involve and attempt to gather evidence in support of a research hypothesis

- In such cases, it is often best to begin with the alternative hypothesis and make it the conclusion that the researcher hopes to support

- The conclusion that the research hypothesis is true is made if the sample data provide sufficient evidence to show that the null hypothesis can be rejected

# Developing Null and Alternative Hypotheses

**Alternative Hypothesis as a Research Hypothesis**

- Example: A new manufacturing method is believed to be better than the current method.

- Alternative Hypothesis:
    - The new manufacturing method is better

- Null Hypothesis:
    - The new methods is no better than the old method

# Developing Null and Alternative Hypotheses

**Alternative Hypothesis as a Research Hypothesis**

- Example: A new bonus plan, that is developed in and attempt to increase sales

- Alternative Hypothesis:
  - The new bonus plan increase sales

- Null Hypothesis:
  - The new bonus plan does not increase sales

# Developing Null and Alternative Hypotheses

**Alternative Hypothesis as a Research Hypothesis**

- Example: A new drug is developed with the goal of lowering Cholesterol-level more than the existing drug

- Alternative Hypothesis:
  - The new drug lowers Cholesterol-level more than the existing drug

- Null Hypothesis:
  - The new drug does not lower Cholesterol-level more than the existing drug

# Developing Null and Alternative Hypotheses

- Null Hypothesis as and assumption to be challenged

- We might begin with a belief or assumption that a statement about the value of a population parameter is true

- **Example**: The label on a milk bottle states that it contains 1000 ml

- Null Hypothesis:
  - The label is correct. $\mu \geq 1000$ ml

- Alternative Hypothesis:
  - The label is incorrect. $\mu < 1000$ ml

# Null and Alternative Hypotheses about a Population Mean μ

- The equality part of the hypotheses always appears in the null hypothesis
- In general, a hypothesis test about the value of a population mean μ must take one of the following three forms (where $\mu_0$ is the hypothesized value of the population mean)

$$H_0: \mu \geq \mu_0$$
$$H_a: \mu < \mu_0$$
One-tailed
(lower-tail)

$$H_0: \mu \leq \mu_0$$
$$H_a: \mu > \mu_0$$
One-tailed
(upper-tail)

$$H_0: \mu = \mu_0$$
$$H_a: \mu \neq \mu_0$$
Two-tailed

# Null and Alternative Hypotheses

- A major hospital in Chennai provides one of the most comprehensive emergency medical services in the world
- Operating in a multiple hospital system with approximately 10 mobile medical units, the service goal is to respond to medical emergencies with a mean time of 8 minutes or less
- The director of medical services wants to formulate a hypothesis test that coils use a sample of emergency response times to determine whether o not the *service goal of 8 minutes or less is being archived.*

# Null and Alternative Hypotheses

The emergency service is meeting the response goal; no follow-up action is necessary.

$$H_0: \quad \mu \leq 8$$

The emergency service is not meeting the response goal; appropriate follow-up action is necessary.

$$H_a: \quad \mu > 8$$

Where: μ = mean response time for the population of medical emergency requests

# Type I Error

- Because hypothesis tests are based on sample data, we must allow for the possibility of errors

- A Type I error is rejecting $H_0$ when it is true

- The probability of making a Type I error when the null hypothesis is called the level of significance

- Applications of hypothesis testing that only control the Type I error are often called significance tests

# Type II Error

- A Type II error is accepting $H_0$ when it is false.

- It is difficult to control for the probability of making a Type II error.

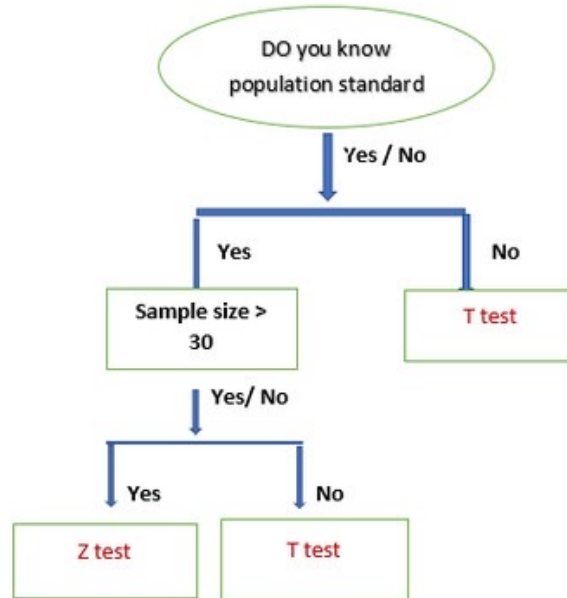- Statisticians avoid the risk of making a Type II error by using "do not reject $H_0$" and not "accept $H_0$".

College Dublin
Computing • IT • Business

# Type I and Type II Errors

| | Population Condition | |
|---|---|---|
| **Conclusion** | **H0 True** $(\mu \leq 8)$ | **H0 False** $(\mu > 8)$ |
| **Accept H0** (Conclude $\mu \leq 8$) | Correct Decision | Type II Error |
| **Reject H0** (Conclude $\mu > 8$) | Type I Error | Correct Decision |

# Three Approached for Hypothesis Testing

- Z test  = Average value

- t- test = = Average value

- Chi-square = categorical data

- ANNOVA = Analysis of variance

# When to use z-test vs t-test

Problem definition 1:

The average height of all players in the academy is 168 cm with a population standard deviation of 39 . New coach believes the mean to be different . He measured height of all 36 players and found average to be 169.5.

A. State Null and Alternate hypothesis

B. At a 95% CI is there enough evidence to accept alternate hypothesis.

# Solution:

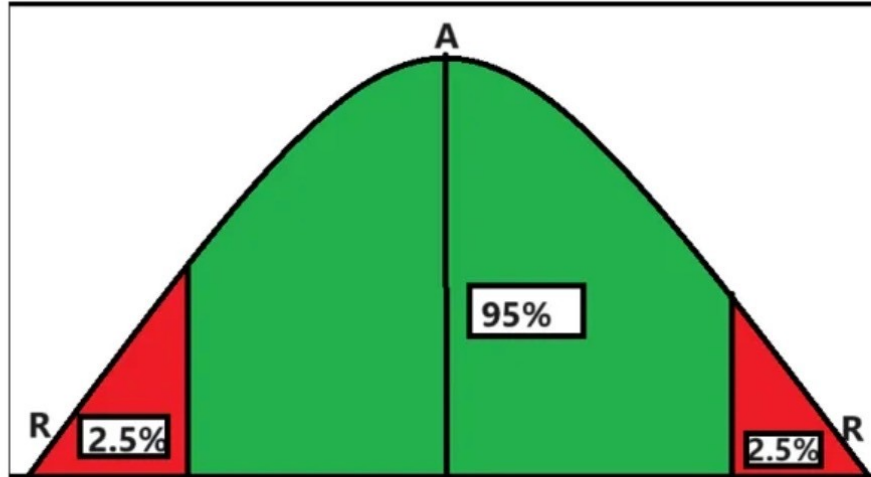$\mu = 168cm$ , $\sigma = 39$ , n= 36 , $\overline{x} = 169.5$

**Null Hypothesis:** H0: $\mu = 168$

**Alternate Hypothesis:** H1: $\mu \neq 168$

Considering above example, it will be two tailed test

**CI = 0.95 => 95%**

$\alpha = 1- CI= 1-0.95 =0.05$

College Dublin
Computing • IT • Business

**As it is 2 tailed it would look like one below**



Here Above A reffers to Accepted Area and R refers to Rejected Area and considering that we have 2 tails 5% is divided in to 2.5% and 2.5%
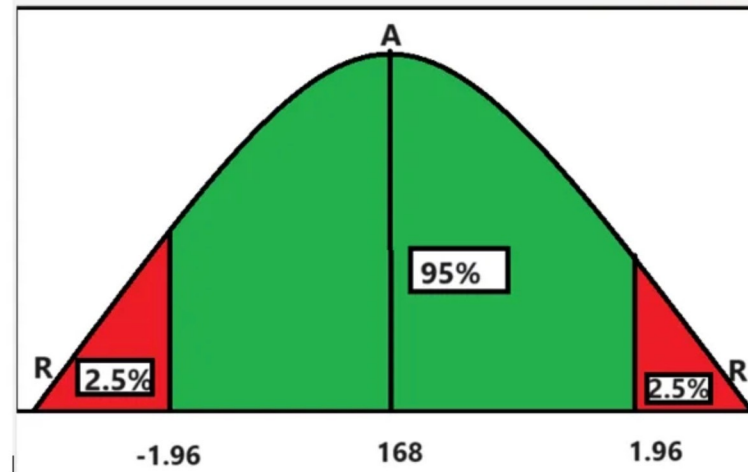
**Decision Boundary:** Now we are calculating decision boundary so we need to find either side as another will inverse of it .

Calculting for right hand curve:

1-0.025= 0.9750 (considering right hand curve so 100%-2.5%)

Refer Z-Table to find value of 0.9750 (Ztable)

Value is 1.96 and inverse will be -1.96



If Z score value falls between -1.96 to 1.96 then we fail to reject Null Hypothesis.

**Z Score:**

$$Z = \frac{(\overline{X} - \mu)}{(\sigma / \sqrt{n})} \qquad SE = \frac{\sigma}{\sqrt{n}}$$

$\mathbf{SE}$ = standard error of the sample

$\sigma$ = sample standard deviation

$n$ = number of samples

$= (169.5 - 169)/(\frac{39}{\sqrt{36}}) \quad = 2.31$

**Conclusion:**

If Z score is greater than 1.96 and lesser then -1.96 then we will reject the null hypothesis.

2.31>1.96 we reject the null hypothesis:

**Hence new coach was right.**

**College Dublin**
Computing • IT • Business

**Problem definition 2**

In the population, the average IQ is 100. A team of scientists wants to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? Use alpha = 0.05.

1. Define Null and Alternative Hypotheses
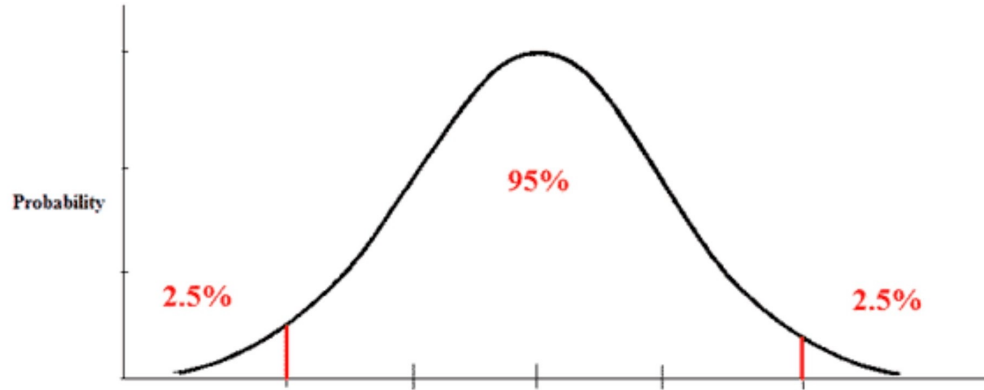
$$H_0; \mu = 100$$

$$H_1; \mu \neq 100$$

**Figure 1.**

2. State Alpha

Alpha = 0.05

3. Calculate Degrees of Freedom

df = n - 1 = 30 - 1 = 29

4. State Decision Rule

Using an alpha of 0.05 with a two-tailed test with 29 degrees of freedom, we would expect our distribution to look something like this:



Use the t-table to look up a two-tailed test with 29 degrees of freedom and an alpha of 0.05. We find a critical value of 2.0452. Thus, our decision rule for this two-tailed test is:

If t is less than -2.0452, or greater than 2.0452, reject the null hypothesis.

5. Calculate Test Statistic

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

$\bar{x} = 140$

$\mu = 100$

$s = 20$

$n = 30$

$$t = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65} = 10.96$$

6. State Results

t = 10.96
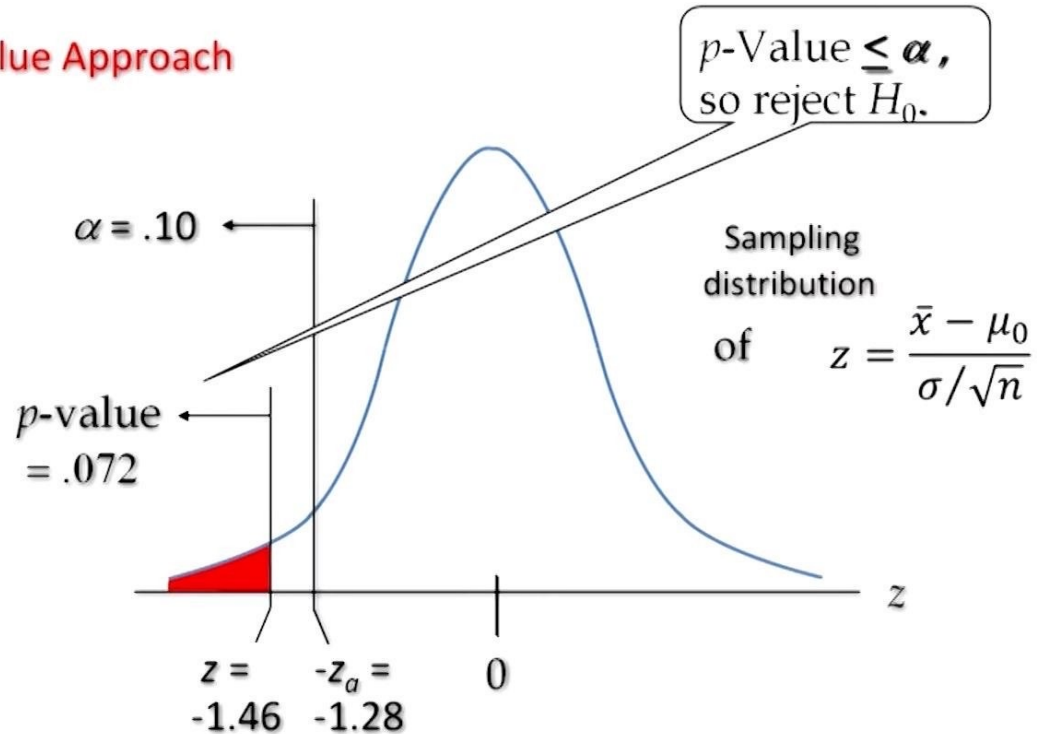
Result: Reject the null hypothesis.

7. State Conclusion

Medication significantly affected intelligence, t = 10.96, p < 0.05.

# p-Value Approach to One-Tailed Hypothesis Testing

- The p-value is the probability, computed using the test statistic, that measures the support (or lack of support) provided by the sample for the null hypothesis

- If the p-value is less than or equal to the level of significance $\alpha$, the value of the test statistic is in the rejection region

- Reject $H_0$ if the p-value ≤ $\alpha$

# Lower-Tailed Test About a Population Mean: σ Known

# p-Value Approach

**Finding P Value**

```
In [3]:  stats.norm.cdf(-1.46)

Out[3]:  0.07214503696589378
```

**Importing library**

```
In [2]:  ▶| from scipy import stats

In [3]:  ▶| stats.norm.cdf(1.96)

Out[3]:  0.9750021048517795
```
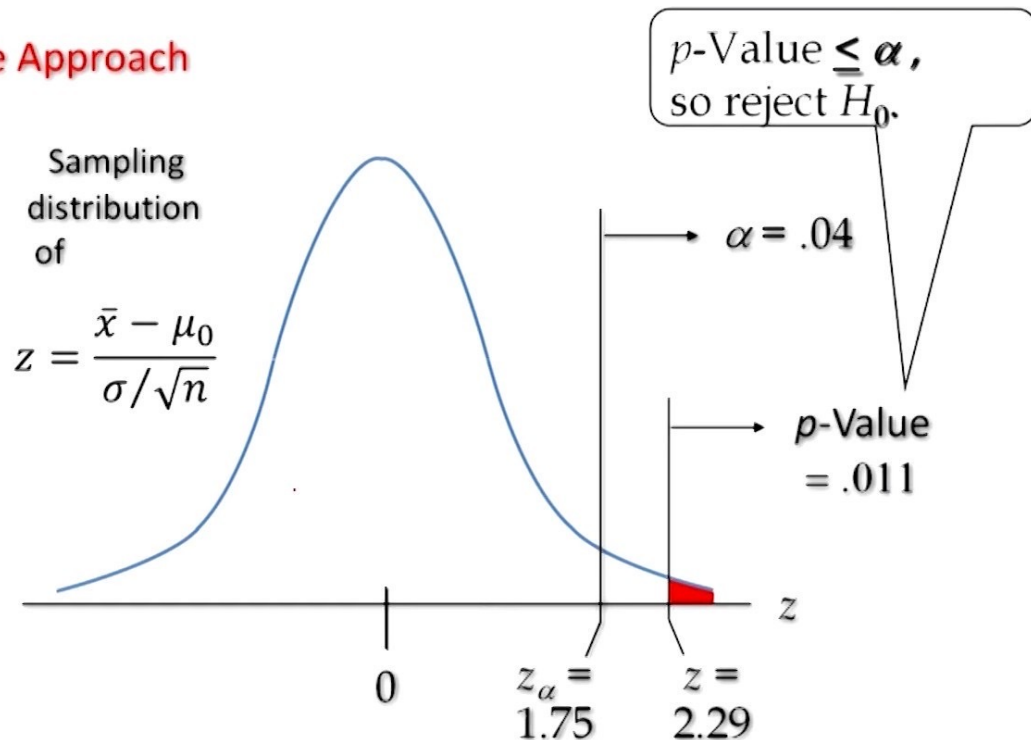
**Finding Z Value**

```
In [5]:  stats.norm.ppf(0.1)

Out[5]:  -1.2815515655446004
```

College Dublin
Computing • IT • Business

# Upper-Tailed Test About a Population Mean: σ Known
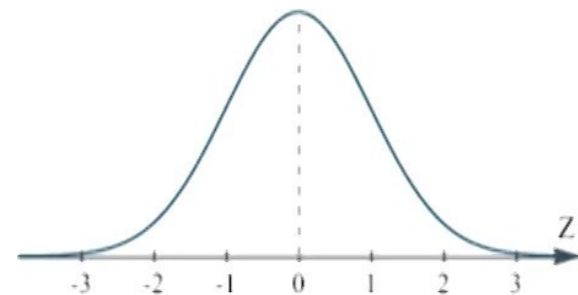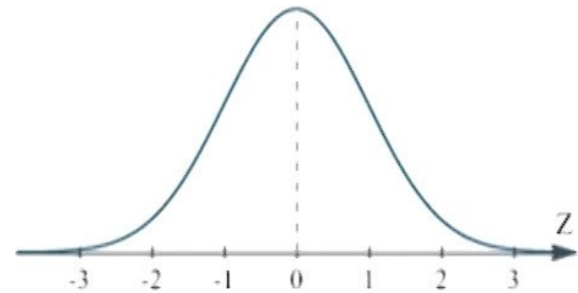
# p-Value Approach

```
In [4]: 1-stats.norm.cdf(1.75)

Out[4]: 0.0400591568638171114
```

```
In [5]: 1-stats.norm.cdf(2.29)

Out[5]: 0.011010658324411393
```
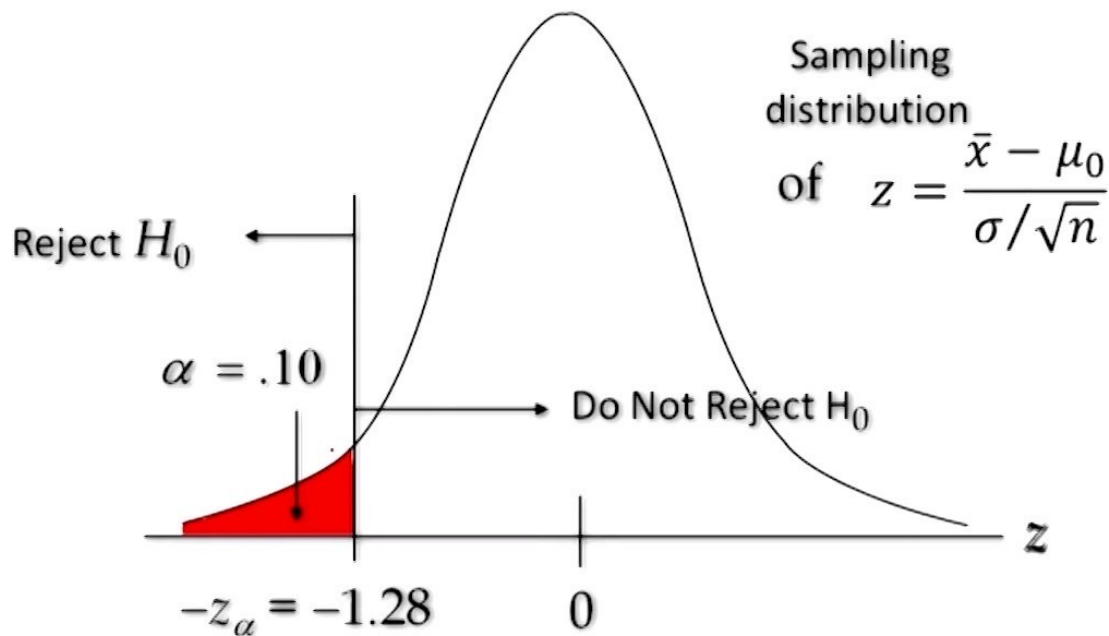
# Critical Value Approach to One-Tailed Hypothesis Testing

- The test statistic z has a standard normal probability distribution.
- We can use the standard normal probability distribution table to find the z-value with and area of α in the lower (or upper) tail of the distribution.

- The value of the test statistic that established the boundary of the rejection region is called the critical value for the test.
- The rejection rule is:
  - Lower tail: Reject $H_0$ if $z \leq -z_\alpha$
  - Upper tail: Reject $H_0$ if $z \geq z_\alpha$

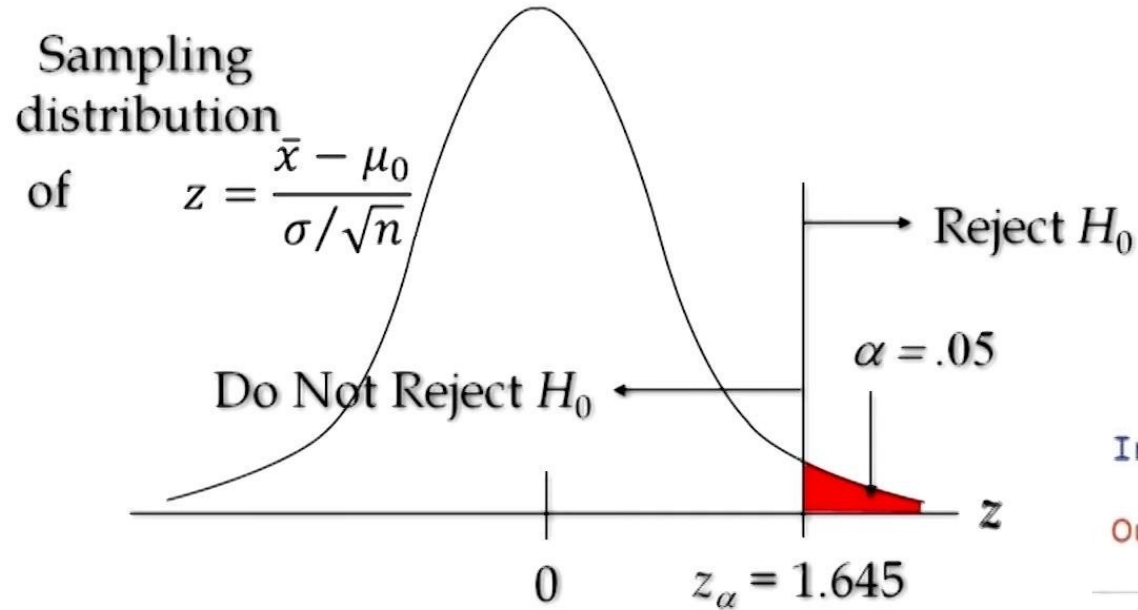# Lower-Tailed Test About a Population Mean: σ

Critical Value Approach



$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

Sampling distribution of

Reject $H_0$

$\alpha = .10$

Do Not Reject $H_0$

$-z_\alpha = -1.28$     $0$

```
In [6]: stats.norm.ppf(0.1)
Out[6]: -1.2815515655446004
```

# Upper-Tailed Test About a Population Mean: σ

Critical Value Approach



Sampling distribution of $z = \dfrac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$

Reject $H_0$

Do Not Reject $H_0$

$\alpha = .05$

$z$

$0$        $z_\alpha = 1.645$

```
In [7]: stats.norm.ppf(0.95)
Out[7]: 1.6448536269514722
```

# Steps of Hypothesis Testing - P value approach

Step 1: Develop the null and alternative hypotheses.

Step 2: Specify the level of significance α.

Step 3: Collect the sample data and compute the test statistic.

P-Value Approach

Step 4: Use the value of the test statistic to compute the p-value.

Step 5: Reject $H_0$ if p-value ≤ α.

# Steps of Hypothesis Testing

Critical Value Approach

Step 4: Use the level of significance α to determine the critical value and the rejection rule.

Step 5: Use the value of the test statistic and the rejection rule to determine whether to reject $H_0$.

# One-Tailed Tests about a population Mean: σ Known

- **Example**: The mean response times for a random sample of 30 Pizza Deliveries is 32 minutes

- The population standard deviation is believed to be 10 minutes.

- The pizza delivery services director wants to perform a hypothesis test, with α = 0.05 level of significance, to determine whether the service goal of 30 minutes or less is begin achieved.

# Given Values

- Sample
- Sample mean = 32 min
- Sample size = 30

- Population
- $\alpha$ = 0.05
- Population mean = 30 min
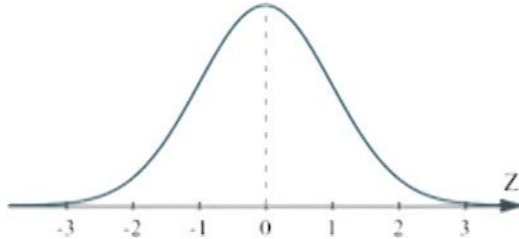
# One-Tailed Tests about a Population Mean: α Known

1.  Develop the hypotheses.
2.  Specify the level of significance.
3.  Compute the value of the test statistic.

$$H_0: \ \mu \leq 30$$
$$H_a: \ \mu > 30$$

$$\alpha = .05$$

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{32 - 30}{10 / \sqrt{30}} = 1.09$$



```
In [8]:  1-stats.norm.cdf(1.09)

Out[8]:  0.1378565720320355
```

# One-Tailed Tests about a Population Mean: α Known
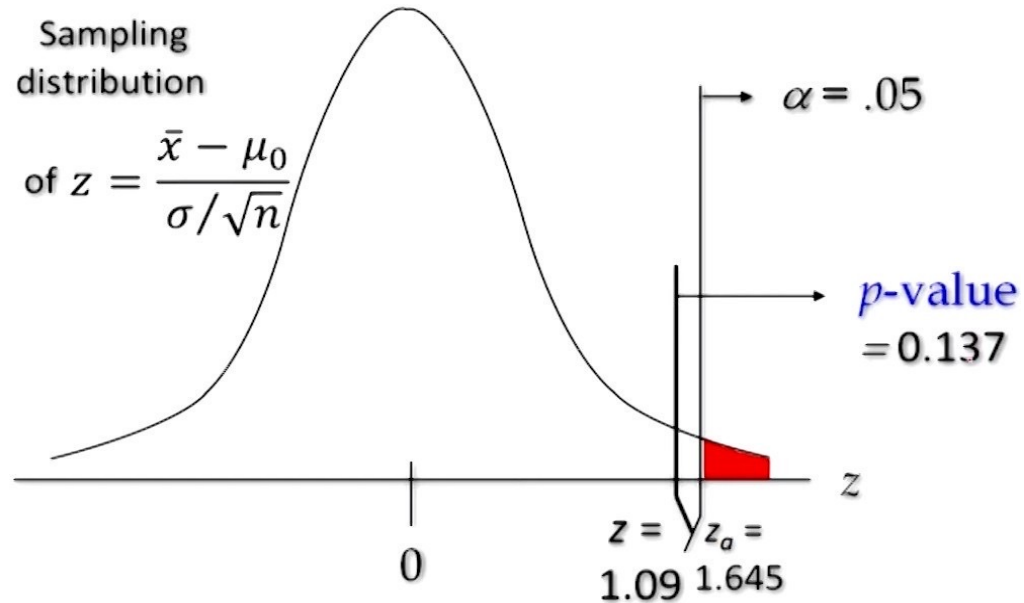## p-Value Approach

4. Compute the p-value,

For z = 1.09, p-value == 0.137

5. Determine whether to reject H0.
- Because p-value = 0.137 > α = .05, we do not reject $H_0$.
- There are not sufficient statistical evidence to infer that Pizza delivery services is not meeting the response goal of 30 minutes.

# One-Tailed Tests about a Population Mean: α Known

Critical Value Approach

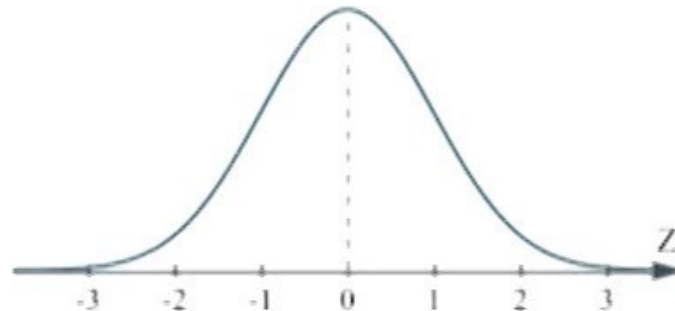# Compute the p-value using the following three steps:

## Critical value Approach

4. Determine the critical value and rejection rule.

For $\alpha$ = .05, $z_{.05}$ = 1.645

Reject $H_0$ if z ≥ 1.645

5. Determine whether to reject $H_0$.

- Because 1.645 ≥ 1.05, we do not reject $H_0$.

# Compute the p-value using the following three steps:

- Compute the value of the test statistic z.
- If z is in the upper tail (z > 0), find the area under the standard normal curve to the right of z.
- If z is in the lower tail (z < 0), find the area under the standard normal curve to the left of z.
- Double the tail area obtained in step 2 to obtain the p -value.
- The rejection rule:
  - Reject $H_0$ of the p-value ≤ α.

# Critical Value Approach to Two-Tailed Hypothesis Testing

- The critical values will occur in both the lower and upper tails of the standard normal curve.

- Use the standard normal probability distribution table to find $z_{\alpha/2}$ (the z-value with and area of $\alpha/2$ in the upper tail of the distribution).

- The rejection rule is:

    Reject $H_0$ if $z \leq -z_{\alpha/2}$ or $z \geq z_{\alpha/2}$.

# Two-Tailed Tests about a Population Mean: σ Known

- **Example**: Milk Carton
- Assume that a sample of 30 milk carton provides a sample mean of 505 ml.
- The population standard deviation is believed to be 10 ml.
- Perform a hypothesis test, at the 0.03 level of significance, population mean 500 ml and help determine whether the filling process should continue operating and corrected.

# Given Values

- Sample
- Sample mean = 505 ml
- Sample size = 30

- Population
- Population mean = 500 ml
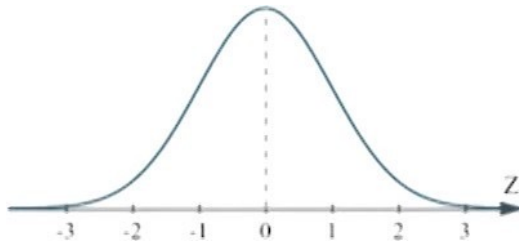- Standard deviation = 10 ml
- Significance level 0.03

# p-value approach

# One-Tailed Tests about a Population Mean: α Known

1. Determine the hypotheses.
2. Specify the level of significance.
3. Compute the value of the test statistic.

$$H_0: \mu = 500$$
$$H_a: \mu \neq 500$$

$$\alpha = .03$$

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{505 - 500}{10 / \sqrt{30}} = 2.74$$



```
In [9]:   1-stats.norm.cdf(2.74)

Out[9]:   0.0030719592186650444


In [10]:  (1-stats.norm.cdf(2.74))*2

Out[10]:  0.006143918437300888
```

# Two-Tailed Tests about a Population Mean: α Known

## p-Value Approach

4. Compute the p-value,

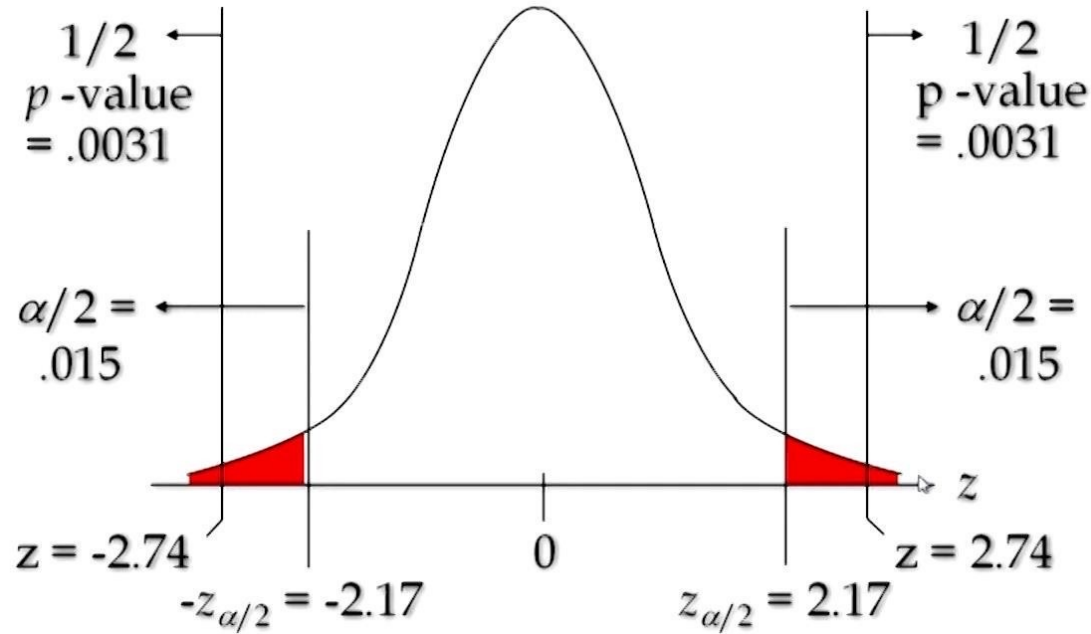$$\text{For } z = 2.74, \text{ p-value} = 2(1 - .9969) = .0061$$

5. Determine whether to reject H0.

- Because p-value = .0062 < α = .03, we reject $H_0$.

There are no sufficient statistical evidence to infer that the null hypothesis is true (i.e. the mean filling quantity is not 500 ml)

# Two-Tailed Tests about a Population Mean: α Known

## p-Value Approach

# Critical Value Approach

# Two-Tailed Tests about a Population Mean: α Known

## Critical value Approach

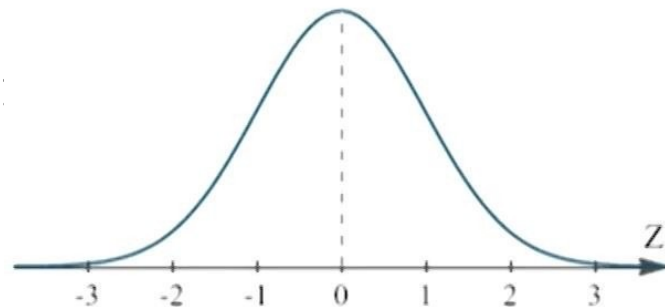4. Determine the critical value and rejection rule.

    For $\alpha/2 = .03/2 = .015$, $z_{0.15} = 2.$

    Reject $H_0$ if $z < -2.17$ or $z > 2.17$

5. Determine whether to reject $H_0$.

● Because $2.74 > 2.17$, we reject $H_0$.

There is sufficient statistical evidence to infer that the null hypothesis is not true



```
In [12]:  stats.norm.ppf(0.015)

Out[12]:  -2.1700903775845606
```
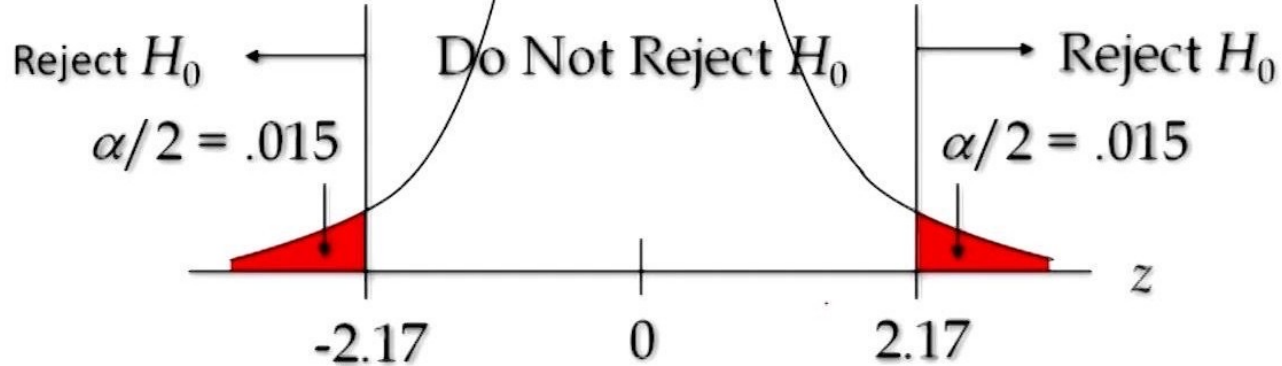
# Two-Tailed Tests about a Population Mean: α Known

**Critical Value Approach**

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{505 - 500}{10 / \sqrt{30}} = 2.74$$

Sampling distribution of $z = \dfrac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$

Reject $H_0$ ← Do Not Reject $H_0$ → Reject $H_0$

$\alpha/2 = .015$

$\alpha/2 = .015$

-2.17    0    2.17    z

Confidence Interval Approach

# Confidence Interval Approach to Two-Tailed Tests about a Population Mean

- Select a simple random sample from the population and use the value of the sample mean to develop the confidence interval for the population mean μ.

- If the confidence interval contains the hypothesis value 500, do not reject $H_0$.

- Otherwise, reject $H_0$.

- Actually, $H_0$ should be rejected if $\mu_0$ happens to be equal to one of the end points of the confidence interval

# Confidence Interval Approach to Two-Tailed Tests about a Population Mean

- The 97% confidence interval for 500 is

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 505 \pm 2.17 \frac{10}{\sqrt{30}} = 505 \pm 3.9619$$
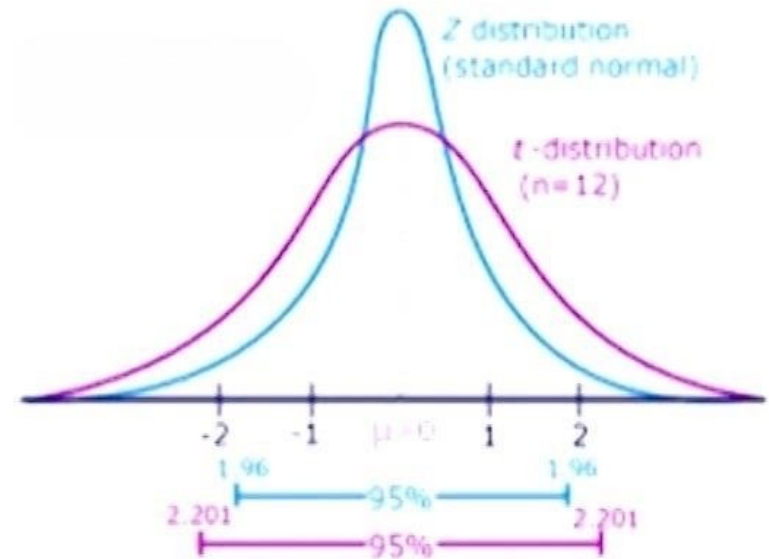$$= \quad 501.03814 , 508.96186$$

- Because the hypothesized value for the population mean, $\mu_0$ = 500 ml, is not in this interval, the hypothesis-testing conclusion is that the null hypothesis, $H_0$: $\mu$ = 500, is rejected.

# Tests About a Population Mean: α Unknown

- Test Statistic

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

- This test statistic has a t distribution with n - 1 degrees of freedom.

# Tests About a Population Mean: α Unknown

- Rejection Rule: p-value Approach

  - Reject H0 if p-value ≤ α

- Rejection Rule: Critical Value Approach

  - H0: $\mu \geq \mu_0$    Reject H0 if $t \leq -t_\alpha$

  - H0: $\mu \leq \mu_0$    Reject H0 if $t \geq t_\alpha$

  - $H_0: \mu = \mu_0$    Reject H0 if $t \leq -t_{\alpha/2}$ or $t \geq t_{\alpha/2}$

```
In [10]:  from scipy import stats
          import numpy as np

In [11]:  x=[10,12,20,21,22,24,18,15]
          stats.ttest_1samp(x,15)

Out[11]:  Ttest_1sampResult(statistic=1.5623450931857947, pvalue=0.1621787560592894)
```

# One-Tailed Test About a Population Mean: α Unknown

**Example: Ice Cream Demand**

- In a ice cream parlor at IIT Roorkee, the following data represent the number of ice-creams sold in 20 days

- Test hypothesis $H_0$: $\mu \leq 10$
- Use $\alpha = .05$ to test the hypothesis.

| Day | No. of Ice-cream Sold | Day | No. of Ice-cream Sold |
|-----|-----|-----|-----|
| 1 | 13 | 11 | 12 |
| 2 | 8 | 12 | 11 |
| 3 | 10 | 13 | 11 |
| 4 | 10 | 14 | 12 |
| 5 | 8 | 15 | 10 |
| 6 | 9 | 16 | 12 |
| 7 | 10 | 17 | 7 |
| 8 | 11 | 18 | 10 |
| 9 | 6 | 19 | 11 |
| 10 | 8 | 20 | 8 |

```
In [8]: x=[13,8,10,10,8,9,10,11,6,8,12,11,11,12,10,12,7,10,11,8]
```
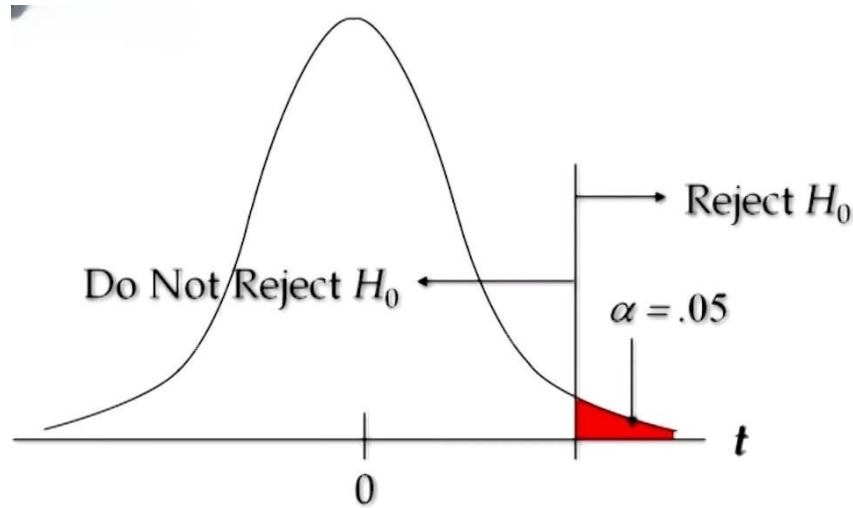
```
In [9]: stats.ttest_1samp(x,10)
```

Out[9]: Ttest_1sampResult(statistic=-0.35843385854878496, pvalue=0.7239703579964252)

```
In [10]: 0.7239703579964252/2
```

Out[10]: 0.3619851789982126

# One-Tailed Test About a Population Mean: α Unknown



```
In [3]:  ▶ stats.t.cdf(-0.384,19)

Out[3]:  0.35262102566795583
```

```
[2]:  ▶ stats.t.ppf(0.05,19)

Out[2]:  -1.7291328115213678
```

Our Dedication.
Your Journey.

cct | College Dublin
Computing • IT • Business