**THOMAS ADEWUNMI UNIVERSITY, OKO, KWARA STATE**
**Faculty of Management and Social Sciences**
**Department of Economics**

**RAIN SEMESTER LECTURE NOTE**        **2023/2024 Session**

**COURSE INFO:**

| | |
|---|---|
| **Course code:** | **ECO 212** |
| **Course title:** | **Applied Statistics II** |
| **Credit unit:** | **2** |

**LECTURER INFO:**

| | |
|---|---|
| **Lecturer's name:** | **Mr Akinbode Damilola** |
| **Department:** | **Economics** |
| **E-mail:** | **damilola.akinbode@tau.edu.ng** |

**TOPIC 1: HYPOTHESIS TESTING**

**Course Description:**

This course note covers the fundamental concepts of hypothesis testing in statistics. It provides an understanding of the steps involved in conducting hypothesis tests, the types of errors, and the interpretation of results. The course also discusses various statistical tests used in hypothesis testing.

**Objectives:**

- Understand the basic concepts and terminology of hypothesis testing.

- Learn the steps involved in conducting a hypothesis test.

- Differentiate between types of errors in hypothesis testing.

- Apply various hypothesis tests to real-world data.

**Course Outline:**

1. **Introduction to Hypothesis Testing**

    o   Definition and Purpose

    o   Null and Alternative Hypotheses

2. **Steps in Hypothesis Testing**

    o   Formulating Hypotheses

    o   Choosing the Significance Level ($\alpha$)

    o   Selecting the Appropriate Test Statistic

    o   Decision Rule and P-value

    o   Drawing Conclusions

3. **Types of Errors in Hypothesis Testing**

    o   Type I Error (False Positive)

    o   Type II Error (False Negative)

    o   Power of a Test

4. **Common Hypothesis Tests**

    o   Z-test

    o   T-test (one-sample, two-sample, paired)

    o   Chi-square Test

    o   ANOVA (Analysis of Variance)

5. **Interpreting Hypothesis Test Results**

    o   P-values and Confidence Intervals

    o   Practical vs. Statistical Significance

6. **Assumptions and Conditions**

    o   Assumptions of Different Tests

    o   Checking for Normality, Independence, and Homogeneity

7. **Case Studies and Applications**

    o   Real-world Examples of Hypothesis Testing

    o   Data Analysis and Interpretation

**Detailed Course Notes:**

# 1. Introduction to Hypothesis Testing

- **Definition and Purpose:**

    - **Hypothesis Testing:** A statistical method used to make decisions or inferences about population parameters based on sample data.

    - Purpose: To determine whether there is enough evidence to reject a null hypothesis in favor of an alternative hypothesis.

- **Null and Alternative Hypotheses:**

    - **Null Hypothesis (H0):** A statement of no effect or no difference, serving as the default assumption.

    - **Alternative Hypothesis (H1 or Ha):** A statement that contradicts the null hypothesis, representing the effect or difference being tested.

# 2. Steps in Hypothesis Testing

- **Formulating Hypotheses:**

    - Define H0 and Ha clearly.

    - Example: H0: $\mu = \mu_0$ (no difference), Ha: $\mu \neq \mu_0$ (difference exists).

- **Choosing the Significance Level ($\alpha$):**

    - Commonly used levels: 0.05, 0.01.

    - Represents the probability of committing a Type I error.

- **Selecting the Appropriate Test Statistic:**

    - Depends on the sample size, data distribution, and type of hypothesis.

    - Examples: Z, t, chi-square, F.

- **Decision Rule and P-value:**

    - **P-value:** The probability of observing the test statistic as extreme as, or more extreme than, the value observed under H0.

    - Compare p-value with $\alpha$: if p-value $\leq \alpha$, reject H0; otherwise, fail to reject H0.

- **Drawing Conclusions:**

    - Based on the comparison of the p-value and the significance level.

    - State the conclusion in the context of the research question.

# 3. Types of Errors in Hypothesis Testing

- **Type I Error (False Positive):**

- o   Rejecting H0 when it is true.

- o   Probability of Type I error = α.

- **Type II Error (False Negative):**

  - o   Failing to reject H0 when it is false.

  - o   Probability of Type II error = β.

- **Power of a Test:**

  - o   The probability of correctly rejecting H0 when it is false (1 - β).

  - o   Higher power is desirable and can be increased by increasing the sample size or effect size.

## 4. Common Hypothesis Tests

- **Z-test:**

  - o   Used for large sample sizes (n > 30) when the population variance is known.

  - o   Example: Testing the population mean.

- **T-test:**

  - o   **One-sample T-test:** Tests if the sample mean is significantly different from a known population mean.

  - o   **Two-sample T-test:** Tests if the means of two independent samples are significantly different.

  - o   **Paired T-test:** Tests if the means of two related groups are significantly different.

- **Chi-square Test:**

  - o   Used for categorical data to test the association between variables or the goodness of fit.

- **ANOVA (Analysis of Variance):**

  - o   Used to compare the means of three or more groups.

  - o   **One-way ANOVA:** Tests for differences among group means.

  - o   **Two-way ANOVA:** Tests for the interaction effect between two independent variables.

## 5. Interpreting Hypothesis Test Results

- **P-values and Confidence Intervals:**

  - o   P-value indicates the strength of evidence against H0.

- Confidence intervals provide a range of values within which the true population parameter is likely to fall.

- **Practical vs. Statistical Significance:**

    - Statistical significance does not always imply practical importance.

    - Consider the effect size and its real-world implications.

## 6. Assumptions and Conditions

- **Assumptions of Different Tests:**

    - Normality: Data should follow a normal distribution for Z and T tests.

    - Independence: Observations should be independent.

    - Homogeneity of Variance: Variances across groups should be equal for ANOVA.

- **Checking for Normality, Independence, and Homogeneity:**

    - Use graphical methods (histograms, Q-Q plots) and statistical tests (Shapiro-Wilk test) to check normality.

    - Check study design and sampling methods for independence.

    - Use Levene's test for homogeneity of variance.

## 7. Case Studies and Applications

- **Real-world Examples of Hypothesis Testing:**

    - Clinical trials comparing the effectiveness of treatments.

    - A/B testing in marketing to compare conversion rates.

- **Data Analysis and Interpretation:**

    - Apply hypothesis tests to datasets using statistical software.

    - Interpret results in the context of the research question and draw actionable insights.

**Recommended Readings:**

- "Statistics for Business and Economics" by Paul Newbold, William L. Carlson, and Betty Thorne

- "Introduction to the Practice of Statistics" by David S. Moore, George P. McCabe, and Bruce A. Craig

- "Statistical Methods for the Social Sciences" by Alan Agresti and Barbara Finlay

# TOPIC 2: Z-TEST HYPOTHESIS TESTING

**Course Description:**

This course note focuses on Z-test hypothesis testing, a statistical method used to determine if there is a significant difference between sample data and a population parameter. It covers the conditions under which Z-tests are applicable, the steps involved in performing a Z-test, and the interpretation of results.

**Objectives:**

- Understand the concept and purpose of Z-tests in hypothesis testing.
- Learn the conditions under which Z-tests are appropriate.
- Apply Z-tests to various hypothesis testing scenarios.
- Interpret the results of Z-tests accurately.

**Course Outline:**

1. **Introduction to Z-Tests**
   o Definition and Purpose
   o Types of Z-Tests

2. **Conditions for Using Z-Tests**
   o Large Sample Sizes
   o Known Population Variance

3. **Steps in Conducting a Z-Test**
   o Formulating Hypotheses
   o Choosing the Significance Level ($\alpha$)
   o Calculating the Z-Statistic
   o Decision Rule and P-value
   o Drawing Conclusions

4. **One-Sample Z-Test**
   o Hypothesis Testing for Population Mean
   o Worked Examples

5. **Two-Sample Z-Test**
   o Hypothesis Testing for the Difference Between Two Means

o   Worked Examples

6. **Interpreting Z-Test Results**

o   P-values and Confidence Intervals

o   Practical vs. Statistical Significance

7. **Common Applications of Z-Tests**

o   Real-world Examples

o   Case Studies

**Detailed Course Notes:**

**1. Introduction to Z-Tests**

- **Definition and Purpose:**

  o   **Z-Test:** A statistical test used to determine whether there is a significant difference between sample data and a population parameter, assuming the data follows a normal distribution.

  o   Purpose: To test hypotheses about population means or proportions using sample data.

- **Types of Z-Tests:**

  o   **One-Sample Z-Test:** Tests if the sample mean is significantly different from a known population mean.

  o   **Two-Sample Z-Test:** Tests if the means of two independent samples are significantly different.

**2. Conditions for Using Z-Tests**

- **Large Sample Sizes:**

  o   Generally, $n > 30$ is considered sufficient for the Central Limit Theorem to apply, allowing the sample mean to be approximately normally distributed.

- **Known Population Variance:**

  o   The population variance ($\sigma^2$) or standard deviation ($\sigma$) must be known. If not, a T-test is more appropriate.

**3. Steps in Conducting a Z-Test**

- **Formulating Hypotheses:**

  o   Define the null hypothesis (H0) and the alternative hypothesis (H1 or Ha).

  o   Example: H0: $\mu = \mu_0$ (no difference), Ha: $\mu \neq \mu_0$ (difference exists).

- **Choosing the Significance Level (α):**
  - Commonly used levels: 0.05, 0.01.
  - Represents the probability of committing a Type I error.
- Calculating the Z-Statistic:
  - Formula for One-Sample Z-Test:

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

  where $\bar{X}$ is the sample mean, $\mu_0$ is the population mean, $\sigma$ is the population standard deviation, and $n$ is the sample size.

  - Formula for Two-Sample Z-Test:

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{(\sigma^2/n_1) + (\sigma^2/n_2)}}$$

  where $\bar{X}_1$ and $\bar{X}_2$ are the sample means, $\mu_1$ and $\mu_2$ are the population means, $\sigma$ is the population standard deviation, and $n_1$ and $n_2$ are the sample sizes.

- **Decision Rule and P-value:**
  - Compare the calculated Z-value to critical Z-values from the standard normal distribution (e.g., ±1.96 for α = 0.05 in a two-tailed test).
  - Alternatively, compare the p-value to the significance level α: if p-value ≤ α, reject H0; otherwise, fail to reject H0.

- **Drawing Conclusions:**
  - Based on the comparison of the p-value and the significance level.
  - State the conclusion in the context of the research question.

## 4. One-Sample Z-Test

- **Hypothesis Testing for Population Mean:**
  - Example Scenario: Testing if the average height of a sample of students is different from the known population mean height.

- **Worked Example:**

  - Population mean (µ0): 65 inches

  - Population standard deviation (σ): 3 inches

  - Sample mean ($\bar{X}$): 66 inches

  - Sample size (n): 50

  - Significance level (α): 0.05

  - Calculate Z-statistic:

$$Z = \frac{66 - 65}{3/\sqrt{50}} \approx 2.36$$

  - Compare Z-value to critical value (±1.96):

    - Since 2.36 > 1.96, reject H0.

  o Conclusion: There is a significant difference between the sample mean and the population mean.

## 5. Two-Sample Z-Test

- **Hypothesis Testing for the Difference Between Two Means:**

  o Example Scenario: Testing if the average scores of two different classes are significantly different.

- Worked Example:

  - Sample mean 1 ($\bar{X}_1$): 78

  - Sample mean 2 ($\bar{X}_2$): 74

  - Population standard deviation ($\sigma$): 10

  - Sample size 1 (n1): 40

  - Sample size 2 (n2): 35

  - Significance level ($\alpha$): 0.05

  - Calculate Z-statistic:

$$Z = \frac{78 - 74}{\sqrt{(10^2/40) + (10^2/35)}} \approx 2.10$$

  - Compare Z-value to critical value ($\pm 1.96$):

    - Since 2.10 > 1.96, reject H0.

  - Conclusion: There is a significant difference between the means of the two classes.

## 6. Interpreting Z-Test Results

- **P-values and Confidence Intervals:**

  - **P-value:** The probability of observing the test statistic as extreme as, or more extreme than, the value observed under H0.

  - **Confidence Interval:** Provides a range of values within which the true population parameter is likely to fall.

- **Practical vs. Statistical Significance:**

  - Statistical significance does not always imply practical importance.

  - Consider the effect size and its real-world implications.

## 7. Common Applications of Z-Tests

- **Real-world Examples:**

  - Comparing the average test scores of a sample to a national average.

- Determining if a new manufacturing process results in a different mean production time compared to the old process.

- **Case Studies:**

  - **Medical Research:** Testing the efficacy of a new drug compared to a standard treatment.

  - **Business Analytics:** Analyzing whether a new marketing strategy significantly impacts sales compared to the previous strategy.

**Recommended Readings:**

- "Statistics for Business and Economics" by Paul Newbold, William L. Carlson, and Betty Thorne

- "Introduction to the Practice of Statistics" by David S. Moore, George P. McCabe, and Bruce A. Craig

- "Fundamentals of Statistics" by Michael Sullivan

# TOPIC 3: T-TEST HYPOTHESIS TESTING

**Course Description:**

This course note provides an in-depth exploration of T-test hypothesis testing, a statistical method used to compare sample means and determine if there are significant differences. It covers various types of T-tests, the conditions under which they are used, and the interpretation of results.

**Objectives:**

- Understand the concept and purpose of T-tests in hypothesis testing.

- Learn the different types of T-tests and their applications.

- Conduct T-tests and interpret their results.

- Apply T-tests to real-world data scenarios.

**Course Outline:**

1. **Introduction to T-Tests**

   o Definition and Purpose

   o Types of T-Tests

2. **Conditions for Using T-Tests**

   o Small Sample Sizes

   o Unknown Population Variance

   o Normality Assumption

3. **Steps in Conducting a T-Test**

   o Formulating Hypotheses

   o Choosing the Significance Level ($\alpha$)

   o Calculating the T-Statistic

   o Decision Rule and P-value

   o Drawing Conclusions

4. **Types of T-Tests**

   o One-Sample T-Test

   o Two-Sample T-Test (Independent Samples)

   o Paired Sample T-Test (Dependent Samples)

5. **One-Sample T-Test**

   o Hypothesis Testing for Population Mean

   o Worked Examples

6. **Two-Sample T-Test**

   o Hypothesis Testing for the Difference Between Two Independent Means

   o Worked Examples

7. **Paired Sample T-Test**

   o Hypothesis Testing for the Difference Between Two Related Means

   o Worked Examples

8. **Interpreting T-Test Results**

   o P-values and Confidence Intervals

   o Practical vs. Statistical Significance

9. **Common Applications of T-Tests**

   o Real-world Examples

   o Case Studies

**Detailed Course Notes:**

**1. Introduction to T-Tests**

- **Definition and Purpose:**

  o **T-Test:** A statistical test used to compare the means of one or two samples and determine if they are significantly different from each other or a known value.

  o Purpose: To test hypotheses about population means using sample data, especially when the population variance is unknown and the sample size is small.

- **Types of T-Tests:**

  o **One-Sample T-Test:** Tests if the sample mean is significantly different from a known population mean.

  o **Two-Sample T-Test:** Tests if the means of two independent samples are significantly different.

  o **Paired Sample T-Test:** Tests if the means of two related groups are significantly different.

**2. Conditions for Using T-Tests**

- **Small Sample Sizes:**

  - T-tests are appropriate for small sample sizes (typically $n < 30$).

- **Unknown Population Variance:**

  - The population variance is unknown, and the sample standard deviation is used as an estimate.

- **Normality Assumption:**

  - The data should approximately follow a normal distribution. For small samples, this is critical; for larger samples, the T-test is robust to deviations from normality.

## 3. Steps in Conducting a T-Test

- **Formulating Hypotheses:**

  - Define the null hypothesis (H0) and the alternative hypothesis (H1 or Ha).

  - Example: H0: $\mu = \mu_0$ (no difference), Ha: $\mu \neq \mu_0$ (difference exists).

- **Choosing the Significance Level (α):**

  - Commonly used levels: 0.05, 0.01.

  - Represents the probability of committing a Type I error.

- **Calculating the T-Statistic:**

  - Formula for One-Sample T-Test:

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

  where $\bar{X}$ is the sample mean, $\mu_0$ is the population mean, $s$ is the sample standard deviation, and $n$ is the sample size.

  - Formula for Two-Sample T-Test (Equal Variances):

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_p\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

  where $s_p$ is the pooled standard deviation, calculated as:

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

  and $n_1$ and $n_2$ are the sample sizes, $s_1$ and $s_2$ are the sample standard deviations.

---

  - Formula for Paired Sample T-Test:

$$t = \frac{\bar{D}}{s_D/\sqrt{n}}$$

  where $\bar{D}$ is the mean of the differences between paired observations, $s_D$ is the standard deviation of the differences, and $n$ is the number of pairs.

- **Decision Rule and P-value:**

  - Compare the calculated t-value to critical t-values from the t-distribution table (based on degrees of freedom and significance level).

  - Alternatively, compare the p-value to the significance level α: if p-value ≤ α, reject H0; otherwise, fail to reject H0.

- **Drawing Conclusions:**

  - Based on the comparison of the p-value and the significance level.

- State the conclusion in the context of the research question.

## 4. Types of T-Tests

### One-Sample T-Test

- **Hypothesis Testing for Population Mean:**
    - Example Scenario: Testing if the average weight of a sample of apples is different from the known population mean weight.

- Worked Example:

    - Population mean ($\mu_0$): 150 grams

    - Sample mean ($\bar{X}$): 155 grams

    - Sample standard deviation (s): 10 grams

    - Sample size (n): 25

    - Significance level ($\alpha$): 0.05

    - Calculate T-statistic:

$$t = \frac{155 - 150}{10/\sqrt{25}} = 2.5$$

    - Degrees of freedom (df): n - 1 = 24

    - Compare t-value to critical value from t-table (±2.064 for df = 24, $\alpha$ = 0.05 two-tailed):

        - Since 2.5 > 2.064, reject H0.

    - Conclusion: There is a significant difference between the sample mean and the population mean.

### Two-Sample T-Test

- **Hypothesis Testing for the Difference Between Two Independent Means:**
    - Example Scenario: Testing if the average heights of male and female students are significantly different.

- **Worked Example:**

  - Sample mean 1 ($\bar{X}_1$): 170 cm

  - Sample mean 2 ($\bar{X}_2$): 165 cm

  - Sample standard deviation 1 (s1): 8 cm

  - Sample standard deviation 2 (s2): 7 cm

  - Sample size 1 (n1): 30

  - Sample size 2 (n2): 25

  - Significance level ($\alpha$): 0.05

  - Calculate pooled standard deviation (sp):

$$s_p = \sqrt{\frac{(30-1)8^2 + (25-1)7^2}{30 + 25 - 2}} \approx 7.54$$

  - Calculate T-statistic:

$$t = \frac{170 - 165}{7.54\sqrt{\frac{1}{30} + \frac{1}{25}}} \approx 2.50$$

  - Degrees of freedom (df): n1 + n2 - 2 = 53

  - Compare t-value to critical value from t-table ($\pm 2.009$ for df = 53, $\alpha = 0.05$ two-tailed):

    - Since 2.50 > 2.009, reject H0.

  - Conclusion: There is a significant difference between the means of the two groups.

**Paired Sample T-Test**

- **Hypothesis Testing for the Difference Between Two Related Means:**

    o Example Scenario: Testing if the average test scores of students before and after a training program are significantly different.

- Worked Example:

    - Mean of differences ($\bar{D}$): 5 points

    - Standard deviation of differences (sD): 3 points

    - Number of pairs (n): 20

    - Significance level (α): 0.05

    - Calculate T-statistic:

$$t = \frac{5}{3/\sqrt{20}} \approx 7.45$$

    - Degrees of freedom (df): n - 1 = 19

    - Compare t-value to critical value from t-table (±2.093 for df = 19, α = 0.05 two-tailed):

        - Since 7.45 > 2.093, reject H0.

    - Conclusion: There is a significant difference between the pre-training and post-training scores.

        ↓

## 8. Interpreting

T-Test Results

- **P-values and Confidence Intervals:**

    o **P-value:** The probability of observing the test statistic as extreme as, or more extreme than, the value observed under H0.

    o **Confidence Interval:** Provides a range of values within which the true population parameter is likely to fall.

- **Practical vs. Statistical Significance:**

    o Statistical significance does not always imply practical importance.

    o Consider the effect size and its real-world implications.

**9. Common Applications of T-Tests**

- **Real-world Examples:**

    o Comparing the effectiveness of two different medications in clinical trials.

    o Analyzing the impact of a new teaching method on student performance.

- **Case Studies:**

    o **Health Sciences:** Testing if there is a difference in blood pressure levels between patients receiving two different treatments.

    o **Education:** Examining if there is a significant difference in reading comprehension scores before and after implementing a new curriculum.

**Recommended Readings:**

- "Statistics for Business and Economics" by Paul Newbold, William L. Carlson, and Betty Thorne

- "Introduction to the Practice of Statistics" by David S. Moore, George P. McCabe, and Bruce A. Craig

- "Fundamentals of Biostatistics" by Bernard Rosner

**TOPIC 4: CHI-SQUARE HYPOTHESIS TESTING**

**Course Description:**

This course note focuses on Chi-Square hypothesis testing, a statistical method used to test relationships between categorical variables. It covers the types of Chi-Square tests, the conditions under which they are appropriate, and the interpretation of results.

**Objectives:**

- Understand the concept and purpose of Chi-Square tests in hypothesis testing.

- Learn the different types of Chi-Square tests and their applications.

- Conduct Chi-Square tests and interpret their results.

- Apply Chi-Square tests to real-world data scenarios.

**Course Outline:**

1. **Introduction to Chi-Square Tests**

   o Definition and Purpose

   o Types of Chi-Square Tests

2. **Conditions for Using Chi-Square Tests**

   o Categorical Data

   o Independence of Observations

   o Expected Frequency Conditions

3. **Steps in Conducting a Chi-Square Test**

   o Formulating Hypotheses

   o Choosing the Significance Level ($\alpha$)

   o Calculating the Chi-Square Statistic

   o Decision Rule and P-value

   o Drawing Conclusions

4. **Types of Chi-Square Tests**

   o Chi-Square Test for Goodness of Fit

   o Chi-Square Test for Independence

   o Chi-Square Test for Homogeneity

5. **Chi-Square Test for Goodness of Fit**

- o   Hypothesis Testing for Distributional Assumptions

- o   Worked Examples

6. **Chi-Square Test for Independence**

- o   Hypothesis Testing for Association Between Two Categorical Variables

- o   Worked Examples

7. **Chi-Square Test for Homogeneity**

- o   Hypothesis Testing for Comparing Proportions Across Groups

- o   Worked Examples

8. **Interpreting Chi-Square Test Results**

- o   P-values and Degrees of Freedom

- o   Practical vs. Statistical Significance

9. **Common Applications of Chi-Square Tests**

- o   Real-world Examples

- o   Case Studies

**Detailed Course Notes:**

**1. Introduction to Chi-Square Tests**

- **Definition and Purpose:**

  - o   **Chi-Square Test:** A statistical test used to determine if there is a significant association between categorical variables or if a sample data matches a population distribution.

  - o   Purpose: To test hypotheses about categorical data.

- **Types of Chi-Square Tests:**

  - o   **Chi-Square Test for Goodness of Fit:** Tests if observed frequencies match expected frequencies based on a specific distribution.

  - o   **Chi-Square Test for Independence:** Tests if there is an association between two categorical variables.

  - o   **Chi-Square Test for Homogeneity:** Tests if different populations have the same proportions of a characteristic.

**2. Conditions for Using Chi-Square Tests**

- **Categorical Data:**

- o   Data should be in categories, not numerical.

- **Independence of Observations:**

  - o   Each observation should be independent of others.

- **Expected Frequency Conditions:**

  - o   Generally, expected frequency in each cell should be at least 5 for the Chi-Square test to be valid.

## 3. Steps in Conducting a Chi-Square Test

- **Formulating Hypotheses:**

  - o   Define the null hypothesis (H0) and the alternative hypothesis (H1 or Ha).

  - o   Example for Goodness of Fit: H0: The data follows a specified distribution, Ha: The data does not follow the specified distribution.

  - o   Example for Independence: H0: There is no association between the variables, Ha: There is an association between the variables.

- **Choosing the Significance Level (α):**

  - o   Commonly used levels: 0.05, 0.01.

  - o   Represents the probability of committing a Type I error.

- Calculating the Chi-Square Statistic:

  - Formula:

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

where $O_i$ is the observed frequency and $E_i$ is the expected frequency.

- **Decision Rule and P-value:**

  - o   Compare the calculated Chi-Square value to critical values from the Chi-Square distribution table (based on degrees of freedom and significance level).

- Alternatively, compare the p-value to the significance level α: if p-value ≤ α, reject H0; otherwise, fail to reject H0.

- **Drawing Conclusions:**

  - Based on the comparison of the p-value and the significance level.

  - State the conclusion in the context of the research question.

## 4. Types of Chi-Square Tests

### Chi-Square Test for Goodness of Fit

- **Hypothesis Testing for Distributional Assumptions:**

  - Example Scenario: Testing if the distribution of colors in a bag of candies matches the manufacturer's claimed distribution.

- Worked Example:

  - Observed frequencies: [50, 30, 20] for colors red, blue, and green.

  - Expected frequencies: [40, 40, 20] based on the manufacturer's claim.

  - Calculate Chi-Square statistic:

$$\chi^2 = \frac{(50-40)^2}{40} + \frac{(30-40)^2}{40} + \frac{(20-20)^2}{20} = 2.5 + 2.5 + 0 = 5.0$$

  - Degrees of freedom (df): Number of categories - 1 = 2

  - Compare Chi-Square value to critical value from Chi-Square table (5.991 for df = 2, α = 0.05):

    - Since 5.0 < 5.991, fail to reject H0.

  - Conclusion: There is no significant difference between the observed and expected distributions.

### Chi-Square Test for Independence

- **Hypothesis Testing for Association Between Two Categorical Variables:**

  - Example Scenario: Testing if gender is independent of voting preference in an election.

- **Worked Example:**
    - Observed frequencies in a 2x2 contingency table:

$$\begin{array}{cc} 30 & 20 \\ 10 & 40 \end{array}$$

    - Expected frequencies:

$$\begin{array}{cc} 20 & 30 \\ 20 & 30 \end{array}$$

    - Calculate Chi-Square statistic:

$$\chi^2 = \frac{(30-20)^2}{20} + \frac{(20-30)^2}{30} + \frac{(10-20)^2}{20} + \frac{(40-30)^2}{30} = 5.0 + 3.33 + 5.0 +$$

    - Degrees of freedom (df): (rows - 1) * (columns - 1) = 1
    - Compare Chi-Square value to critical value from Chi-Square table (3.841 for df = 1, α = 0.05): ↓

  - Since 16.66 > 3.841, reject H0.

  o Conclusion: There is a significant association between gender and voting preference.

**Chi-Square Test for Homogeneity**

- **Hypothesis Testing for Comparing Proportions Across Groups:**
  o Example Scenario: Testing if different schools have the same proportion of students passing a standardized test.

- **Worked Example:**
    - Observed frequencies:

$$\begin{array}{cc} 40 & 10 \\ 50 & 10 \\ 30 & 20 \end{array}$$

    - Calculate expected frequencies and Chi-Square statistic similarly to the test for independence.
    - Degrees of freedom (df): (number of groups - 1) * (number of categories - 1). ↓

**8. Interpreting Chi-Square Test Results**

- **P-values and Degrees of Freedom:**

  - **P-value:** The probability of observing the test statistic as extreme as, or more extreme than, the value observed under H0.

  - **Degrees of Freedom:** Determines the critical value from the Chi-Square distribution table.

- **Practical vs. Statistical Significance:**

  - Statistical significance does not always imply practical importance.

  - Consider the effect size and its real-world implications.

**9. Common Applications of Chi-Square Tests**

- **Real-world Examples:**

  - Testing if customer satisfaction is independent of store location.

  - Analyzing if the distribution of a genetic trait matches expected Mendelian ratios.

- **Case Studies:**

  - **Market Research:** Testing if purchasing preferences are independent of demographic factors.

  - **Healthcare:** Examining if the occurrence of a disease is independent of geographic region.

**Recommended Readings:**

- "Statistics for Business and Economics" by Paul Newbold, William L. Carlson, and Betty Thorne

- "Introduction to the Practice of Statistics" by David S. Moore, George P. McCabe, and Bruce A. Craig

- "Fundamentals of Biostatistics" by Bernard Rosner

# TOPIC 5: CORRELATION ANALYSIS

**Course Description:**

This course note provides a comprehensive understanding of correlation analysis, a statistical method used to measure the strength and direction of the relationship between two quantitative variables. It covers the types of correlation coefficients, methods for calculating them, and interpretation of the results.

**Objectives:**

- Understand the concept and purpose of correlation analysis.

- Learn the different types of correlation coefficients and their applications.

- Calculate correlation coefficients using different methods.

- Interpret the results of correlation analysis.

- Apply correlation analysis to real-world data scenarios.

**Course Outline:**

1. **Introduction to Correlation Analysis**

   o Definition and Purpose

   o Types of Correlation

2. **Types of Correlation Coefficients**

   o Pearson's Correlation Coefficient

   o Spearman's Rank Correlation Coefficient

   o Kendall's Tau Correlation Coefficient

3. **Conditions for Using Correlation Analysis**

   o Linearity

   o Level of Measurement

   o Homoscedasticity

4. **Steps in Conducting Correlation Analysis**

   o Data Collection

   o Checking Assumptions

   o Calculating the Correlation Coefficient

   o Hypothesis Testing

o　Interpreting Results

5. **Pearson's Correlation Coefficient**

　　　o　Definition and Formula

　　　o　Calculation and Interpretation

　　　o　Worked Examples

6. **Spearman's Rank Correlation Coefficient**

　　　o　Definition and Formula

　　　o　Calculation and Interpretation

　　　o　Worked Examples

7. **Kendall's Tau Correlation Coefficient**

　　　o　Definition and Formula

　　　o　Calculation and Interpretation

　　　o　Worked Examples

8. **Interpreting Correlation Coefficients**

　　　o　Strength and Direction

　　　o　Significance Testing

　　　o　Practical vs. Statistical Significance

9. **Common Applications of Correlation Analysis**

　　　o　Real-world Examples

　　　o　Case Studies

**Detailed Course Notes:**

**1. Introduction to Correlation Analysis**

- **Definition and Purpose:**

　　　o　**Correlation Analysis:** A statistical technique used to quantify the relationship between two variables.

　　　o　Purpose: To determine the strength and direction of the relationship between variables and to identify patterns.

- **Types of Correlation:**

　　　o　**Positive Correlation:** Both variables move in the same direction.

- o **Negative Correlation:** Variables move in opposite directions.
- o **No Correlation:** No discernible relationship between variables.

## 2. Types of Correlation Coefficients

- **Pearson's Correlation Coefficient (r):**
  - o Measures linear relationships between two continuous variables.
  - o Values range from -1 to 1, where -1 indicates a perfect negative linear relationship, 1 indicates a perfect positive linear relationship, and 0 indicates no linear relationship.

- **Spearman's Rank Correlation Coefficient ($\rho$ or rs):**
  - o Measures the strength and direction of the relationship between two ranked variables.
  - o Suitable for ordinal data or non-linear relationships.

- **Kendall's Tau Correlation Coefficient ($\tau$):**
  - o Measures the association between two variables based on the ranks of data.
  - o Less sensitive to errors in data than Spearman's coefficient.

## 3. Conditions for Using Correlation Analysis

- **Linearity:**
  - o Pearson's correlation requires a linear relationship between variables.
  - o Spearman's and Kendall's correlations do not require linearity.

- **Level of Measurement:**
  - o Pearson's: Interval or ratio data.
  - o Spearman's and Kendall's: Ordinal, interval, or ratio data.

- **Homoscedasticity:**
  - o The spread of one variable is consistent across the range of another variable.

## 4. Steps in Conducting Correlation Analysis

- **Data Collection:**
  - o Gather data for the variables of interest.

- **Checking Assumptions:**
  - o Ensure the data meets the assumptions for the chosen correlation coefficient.

- **Calculating the Correlation Coefficient:**

  o Use appropriate formulas or statistical software.

- **Hypothesis Testing:**

  o Null hypothesis (H0): There is no correlation between the variables.

  o Alternative hypothesis (H1): There is a correlation between the variables.

- **Interpreting Results:**

  o Determine the strength and direction of the relationship.

  o Assess statistical significance.

## 5. Pearson's Correlation Coefficient

- **Definition and Formula:**

  o Measures the linear relationship between two continuous variables.

- Formula:

$$r = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum(X_i - \bar{X})^2 \sum(Y_i - \bar{Y})^2}}$$

where $X_i$ and $Y_i$ are the values of the variables, and $\bar{X}$ and $\bar{Y}$ are the means of the variables.

- Calculation and Interpretation:

  - Calculate using the formula or statistical software.

  - Interpret the value of $r$ to determine the strength and direction of the relationship.

    - $0.0 - 0.3$: Weak correlation

    - $0.3 - 0.7$: Moderate correlation

    - $0.7 - 1.0$: Strong correlation

- **Worked Examples:**

  o Example Scenario: Analyzing the correlation between study hours and exam scores.

  o Given data for study hours and scores, calculate r and interpret the result.

## 6. Spearman's Rank Correlation Coefficient

- **Definition and Formula:**

  - Measures the strength and direction of the relationship between two ranked variables.

- Formula:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

where $d_i$ is the difference between the ranks of corresponding values and $n$ is the number of observations.

Calculation and Interpretation:

- Rank the data, calculate the differences, and use the formula.

- Interpret the value of $\rho$ to determine the strength and direction of the relationship.

- Formula:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

where $d_i$ is the difference between the ranks of corresponding values and $n$ is the number of observations.

Calculation and Interpretation:

- Rank the data, calculate the differences, and use the formula.

- Interpret the value of $\rho$ to determine the strength and direction of the relationship.

- **Worked Examples:**

  - Example Scenario: Analyzing the correlation between employee ranks and performance scores.

  - Given ranked data, calculate $\rho$ and interpret the result.

## 7. Kendall's Tau Correlation Coefficient

- **Definition and Formula:**

o   Measures the association between two variables based on the ranks of data.

- Formula:

$$\tau = \frac{(C - D)}{\frac{1}{2}n(n - 1)}$$

where $C$ is the number of concordant pairs and $D$ is the number of discordant pairs.

Calculation and Interpretation:

- Determine concordant and discordant pairs, and use the formula.

- Interpret the value of $\tau$ to determine the strength and direction of the relationship.

- **Worked Examples:**

    o   Example Scenario: Analyzing the correlation between customer satisfaction ranks and loyalty ranks.

    o   Given ranked data, calculate $\tau$ and interpret the result.

## 8. Interpreting Correlation Coefficients

- **Strength and Direction:**

    o   Positive values indicate a positive relationship; negative values indicate a negative relationship.

    o   The closer the value is to $\pm 1$, the stronger the relationship.

- **Significance Testing:**

    o   Use p-values to determine if the correlation is statistically significant.

    o   Null hypothesis: There is no correlation between the variables.

    o   If p-value $\leq \alpha$, reject the null hypothesis.

- **Practical vs. Statistical Significance:**

    o   Statistical significance does not always imply practical importance.

    o   Consider the effect size and its real-world implications.

## 9. Common Applications of Correlation Analysis

- **Real-world Examples:**

- Examining the relationship between marketing spend and sales revenue.

- Analyzing the correlation between temperature and ice cream sales.

- **Case Studies:**

  - **Finance:** Studying the correlation between stock prices of two companies.

  - **Healthcare:** Investigating the correlation between BMI and blood pressure levels.

**Recommended Readings:**

- "Statistics for Business and Economics" by Paul Newbold, William L. Carlson, and Betty Thorne

- "Introduction to the Practice of Statistics" by David S. Moore, George P. McCabe, and Bruce A. Craig

- "Fundamentals of Biostatistics" by Bernard Rosner

# TOPIC 6: REGRESSION ANALYSIS

**Course Description:**

This course note provides a comprehensive understanding of regression analysis, a statistical method used to examine the relationship between a dependent variable and one or more independent variables. It covers simple linear regression, multiple linear regression, assumptions, interpretation, and applications.

**Objectives:**

- Understand the concept and purpose of regression analysis.

- Learn the different types of regression models and their applications.

- Conduct regression analysis using statistical software.

- Interpret the results of regression analysis.

- Apply regression analysis to real-world data scenarios.

**Course Outline:**

1. **Introduction to Regression Analysis**
   - Definition and Purpose
   - Types of Regression Models

2. **Simple Linear Regression**
   - Model Specification
   - Estimation of Parameters
   - Assumptions of Simple Linear Regression

3. **Multiple Linear Regression**
   - Model Specification
   - Estimation of Parameters
   - Assumptions of Multiple Linear Regression

4. **Assumptions of Regression Analysis**
   - Linearity
   - Independence
   - Homoscedasticity
   - Normality

5. **Steps in Conducting Regression Analysis**

   o Data Collection

   o Model Specification

   o Estimation of Parameters

   o Hypothesis Testing

   o Model Diagnostics and Validation

6. **Interpreting Regression Results**

   o Coefficients and Significance

   o R-squared and Adjusted R-squared

   o Residual Analysis

   o Practical vs. Statistical Significance

7. **Common Applications of Regression Analysis**

   o Real-world Examples

   o Case Studies

**Detailed Course Notes:**

**1. Introduction to Regression Analysis**

- **Definition and Purpose:**

  o **Regression Analysis:** A statistical technique used to examine the relationship between a dependent variable and one or more independent variables.

  o Purpose: To predict the value of the dependent variable based on the values of the independent variables and to understand the relationships between variables.

- **Types of Regression Models:**

  o **Simple Linear Regression:** One independent variable.

  o **Multiple Linear Regression:** More than one independent variable.

  o **Logistic Regression:** Dependent variable is categorical.

  o **Polynomial Regression:** Non-linear relationship between variables.

**2. Simple Linear Regression**

- **Model Specification:**
    - Linear relationship between a single independent variable (X) and a dependent variable (Y).
    - Equation: $Y = \beta_0 + \beta_1 X + \epsilon$
        - $Y$: Dependent variable
        - $X$: Independent variable
        - $\beta_0$: Intercept
        - $\beta_1$: Slope
        - $\epsilon$: Error term

- **Estimation of Parameters:**
    - **Ordinary Least Squares (OLS) Method:** Minimizes the sum of squared residuals to estimate $\beta_0$ and $\beta_1$.
    - Formulas:

$$\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

- **Assumptions of Simple Linear Regression:**
    - Linearity: The relationship between X and Y is linear.
    - Independence: Observations are independent.
    - Homoscedasticity: Constant variance of errors.
    - Normality: Errors are normally distributed.

## 3. Multiple Linear Regression

- **Model Specification:**

  - Linear relationship between multiple independent variables (X1, X2, ..., Xp) and a dependent variable (Y).

  - Equation: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + ... + \beta_p X_p + \epsilon$

- **Estimation of Parameters:**

  - **Ordinary Least Squares (OLS) Method:** Minimizes the sum of squared residuals to estimate the parameters.

  - Matrix notation:

$$\hat{\beta} = (X'X)^{-1}X'Y$$

- **Assumptions of Multiple Linear Regression:**

  o Linearity: The relationship between each X and Y is linear.

  o Independence: Observations are independent.

  o Homoscedasticity: Constant variance of errors.

  o Normality: Errors are normally distributed.

  o No multicollinearity: Independent variables are not highly correlated.

## 4. Assumptions of Regression Analysis

- **Linearity:** The relationship between the independent and dependent variables is linear.

- **Independence:** Observations of the dependent variable are independent of each other.

- **Homoscedasticity:** The variance of the error terms is constant across all levels of the independent variables.

- **Normality:** The error terms are normally distributed.

## 5. Steps in Conducting Regression Analysis

- **Data Collection:**

  o Gather data for the dependent and independent variables.

- **Model Specification:**

  o Define the regression model to be used based on the data and research question.

- **Estimation of Parameters:**

  o Use statistical software to estimate the parameters of the regression model.

- **Hypothesis Testing:**

  o Null hypothesis (H0): No relationship between the independent and dependent variables.

  o Alternative hypothesis (H1): There is a relationship between the independent and dependent variables.

  o Use t-tests for individual coefficients and F-tests for the overall model.

- **Model Diagnostics and Validation:**

  o Check for violations of regression assumptions.

  o Use diagnostic plots (e.g., residual plots) and statistical tests (e.g., Breusch-Pagan test for homoscedasticity).

## 6. Interpreting Regression Results

- **Coefficients and Significance:**

  o Coefficients ($\beta$): Indicate the direction and strength of the relationship between each independent variable and the dependent variable.

  o p-values: Assess the statistical significance of each coefficient.

- **R-squared and Adjusted R-squared:**

  o **R-squared ($R^2$):** Proportion of variance in the dependent variable explained by the independent variables.

  o **Adjusted R-squared:** Adjusts $R^2$ for the number of predictors in the model.

- **Residual Analysis:**

  o Analyze residuals to check for patterns that might indicate violations of regression assumptions.

- **Practical vs. Statistical Significance:**

  o Statistical significance does not always imply practical importance.

  o Consider the effect size and its real-world implications.

## 7. Common Applications of Regression Analysis

- **Real-world Examples:**

  o Predicting house prices based on features like size, location, and number of bedrooms.

  o Analyzing the impact of advertising spend on sales revenue.

- **Case Studies:**

- o **Economics:** Studying the relationship between GDP growth and unemployment rates.

- o **Healthcare:** Investigating the impact of lifestyle factors on health outcomes.

**Recommended Readings:**

- "Applied Regression Analysis" by Norman R. Draper and Harry Smith

- "Regression Analysis by Example" by Samprit Chatterjee and Ali S. Hadi

- "An Introduction to Statistical Learning" by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani