

inferential-statistics-assignment-jun6

July 6, 2024

1 INFERENTIAL STATISTICS

[1]: *#Inferential Statistics*

2 INTRODUCTION

Inferential statistics is a branch of statistics that allows us to make conclusions or inferences about a population based on data from a sample of that population. It involves the use of various techniques and methods to analyze sample data, make predictions, and test hypotheses. Here are the key components and concepts of inferential statistics:

1. Population vs. Sample Population: The entire group of individuals or instances about whom we hope to learn. Sample: A subset of the population, selected for measurement and analysis.
2. Parameters and Statistics Parameter: A numerical summary of a population (e.g., population mean μ , population proportion P). Statistic: A numerical summary of a sample (e.g., sample mean \bar{x} , sample proportion \hat{P}).
3. Sampling Methods Random Sampling: Each member of the population has an equal chance of being selected. Stratified Sampling: The population is divided into strata, and random samples are taken from each stratum. Cluster Sampling: The population is divided into clusters, and entire clusters are randomly selected. Systematic Sampling: Every k th member of the population is selected.
4. Hypothesis Testing Null Hypothesis (H_0): A statement of no effect or no difference, which we test against. Alternative Hypothesis (H_1): A statement that indicates the presence of an effect or a difference. P-value: The probability of observing the sample data, or something more extreme, assuming the null hypothesis is true. Significance Level (α): The threshold at which we reject the null hypothesis (commonly 0.05).
5. Confidence Intervals A range of values, derived from the sample statistics, that is likely to contain the population parameter. Formula: $\bar{x} \pm z \frac{s}{\sqrt{n}}$

where \bar{x} is the sample mean, z is the z-score corresponding to the desired confidence level, s is the sample standard deviation, and n is the sample size. 6. Common Tests and Techniques Z-Test: Used for hypothesis testing when the population variance is known and the sample size is large. T-Test: Used when the population variance is unknown and the sample size is small. One-sample T-test: Tests whether the sample mean is different from a known value. Two-sample T-test: Tests whether the means of two independent samples are different. Paired T-test: Tests whether

the means of two related samples are different. ANOVA (Analysis of Variance): Used to compare the means of three or more samples. Chi-Square Test: Used for categorical data to assess how likely it is that an observed distribution is due to chance. Regression Analysis: Used to examine the relationship between variables. 7. Assumptions in Inferential Statistics Random Sampling: The sample should be randomly selected from the population. Independence: Observations must be independent of each other. Normality: Data should be approximately normally distributed (especially for small sample sizes). Homogeneity of Variance: Variances within different groups should be roughly equal. 8. Errors in Hypothesis Testing Type I Error (α): Rejecting the null hypothesis when it is true (false positive). Type II Error (β): Failing to reject the null hypothesis when it is false (false negative). Applications of Inferential Statistics Inferential statistics is widely used in various fields such as:

Medicine: For clinical trials and health studies to determine the effectiveness of treatments.

Economics:

To forecast economic trends and make policy decisions. Psychology:

To understand behavior patterns and psychological phenomena. Marketing: To analyze consumer behavior and improve marketing strategies.

Conclusion

Inferential statistics provides the tools to make informed decisions based on sample data. By understanding the principles and methods of inferential statistics, researchers can draw meaningful conclusions about a population, despite only having a subset of data.

3 Inferential Statistics with Python

1. Introduction to Parameters and Statistics Parameter: A measure (mean, median, variance, etc.) for population data. Statistic: A measure (mean, median, variance, etc.) for sample data.

In practice, calculating parameters for entire populations is often infeasible due to size and cost constraints, so we use sample statistics to estimate population parameters.

Example: Estimating the average life expectancy in India using a sample. 2. Sampling Techniques

a) Simple Random Sampling

In simple random sampling, every member of the population has an equal chance of being selected.

Example: Selecting 12 tomato plants from a farm of 98 for a growth study.

```
[2]: import numpy as np

# Total population
population = np.arange(1, 99)
# Sample size
sample_size = 12

# Simple random sampling without replacement
```

```

sample_without_replacement = np.random.choice(population, sample_size,
↪replace=False)
print("Sample without replacement:", sample_without_replacement)

# Simple random sampling with replacement
sample_with_replacement = np.random.choice(population, sample_size,
↪replace=True)
print("Sample with replacement:", sample_with_replacement)

```

Sample without replacement: [5 84 93 89 79 3 11 87 45 40 17 22]
Sample with replacement: [46 19 82 41 10 1 89 2 83 93 4 61]

4 b) Stratified Sampling

In stratified sampling, the population is divided into strata, and random samples are taken from each stratum.

Example: Selecting 2 roses of each color from a nursery.

```

[3]: import pandas as pd

# Example data
roses = ['White', 'Pink', 'White', 'Red', 'Yellow', 'Orange', 'Orange', 'Red',
↪'Yellow', 'White', 'Pink', 'White', 'Red', 'Orange']
roses_df = pd.DataFrame({'Color': roses})

# Grouping by color
grouped = roses_df.groupby('Color').apply(lambda x: x.sample(2)).
↪reset_index(drop=True)
print(grouped)

```

```

      Color
0  Orange
1  Orange
2    Pink
3    Pink
4     Red
5     Red
6   White
7   White
8  Yellow
9  Yellow

```

5 c) Systematic Sampling

In systematic sampling, the first data point is selected randomly, and subsequent points are selected at regular intervals.

Example: Ann selecting 4 marbles starting from the 2nd.

```
[4]: # Total population
marbles = np.arange(1, 21)
# Start point
start_point = 2
# Interval
interval = 5
# Sample size
sample_size = 4

# Systematic sampling
sample_systematic = marbles[start_point-1:start_point-1+sample_size*interval:
    ↪interval]
print("Systematic sample:", sample_systematic)
```

Systematic sample: [2 7 12 17]

6 3. Central Limit Theorem (CLT)

For a sufficiently large sample size (typically $n > 30$), the distribution of sample means approximates a normal distribution, regardless of the population's distribution.

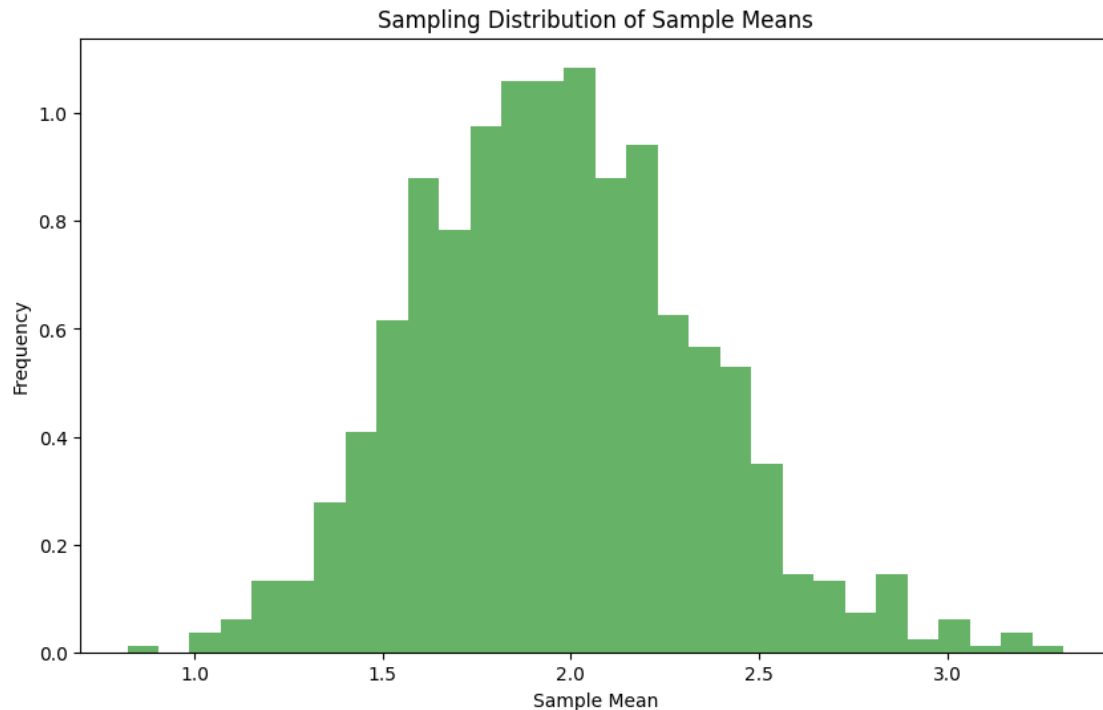
Example: Demonstrating CLT with Python.

```
[5]: import matplotlib.pyplot as plt

# Simulate population data
population = np.random.exponential(scale=2, size=1000)

# Draw samples and compute means
sample_means = [np.mean(np.random.choice(population, 30)) for _ in range(1000)]

# Plotting
plt.figure(figsize=(10, 6))
plt.hist(sample_means, bins=30, density=True, alpha=0.6, color='g')
plt.title('Sampling Distribution of Sample Means')
plt.xlabel('Sample Mean')
plt.ylabel('Frequency')
plt.show()
```



7 4. Estimation

a) Point Estimation

Using sample statistics to estimate population parameters.

Example: Estimating the mean length of steel rods.

```
[6]: # Sample data
length_rod = [25.2, 26.3, 28, 21.9, 23.4, 24, 27.2, 23, 29.2, 28.7, 23.1, 23.5, 26.4, 22.8, 24.7]
sample_mean = np.mean(length_rod)
print("Sample mean:", sample_mean)
```

Sample mean: 25.16

8 b) Interval Estimation (Confidence Intervals)

Estimating a range within which the population parameter likely falls.

Example: Calculating a 90% confidence interval for average sugar content in baby food.

```
[7]: import scipy.stats as stats

# Given data
sample_mean = 24
std_dev = 8
sample_size = 37
confidence_level = 0.90

# Z-score for 90% confidence level
z_score = stats.norm.isf((1 - confidence_level) / 2)
margin_of_error = z_score * (std_dev / np.sqrt(sample_size))

# Confidence interval
confidence_interval = (sample_mean - margin_of_error, sample_mean +
    ↪margin_of_error)
print("90% Confidence Interval:", confidence_interval)
```

90% Confidence Interval: (21.83670183570907, 26.16329816429093)

9 5. Hypothesis Testing

- a) Null Hypothesis (H_0): There is no effect or difference.
- b) Alternative Hypothesis (H_a): There is an effect or difference.

Example: Testing the weight of green tea bags.

```
[8]: # Given data
sample_mean = 3.28
population_mean = 3.5
std_dev = 0.6
sample_size = 50
alpha = 0.10

# Test statistic (Z-score)
z_score = (sample_mean - population_mean) / (std_dev / np.sqrt(sample_size))

# P-value
p_value = stats.norm.sf(z_score)
print("P-value:", p_value)

# Conclusion
if p_value < alpha:
    print("Reject the null hypothesis")
else:
    print("Fail to reject the null hypothesis")
```

P-value: 0.9952390544079506
Fail to reject the null hypothesis

10 Summary

Sampling Techniques: Simple random, stratified, and systematic sampling methods help in selecting representative samples. Central Limit Theorem: Demonstrates that the distribution of sample means is approximately normal for large samples. Estimation: Point estimation and confidence intervals are used to infer population parameters. Hypothesis Testing: Involves testing claims about population parameters using sample data. These concepts are fundamental to inferential statistics and form the basis for more advanced statistical techniques and data analysis in various fields.

11 Graphs and Visualization

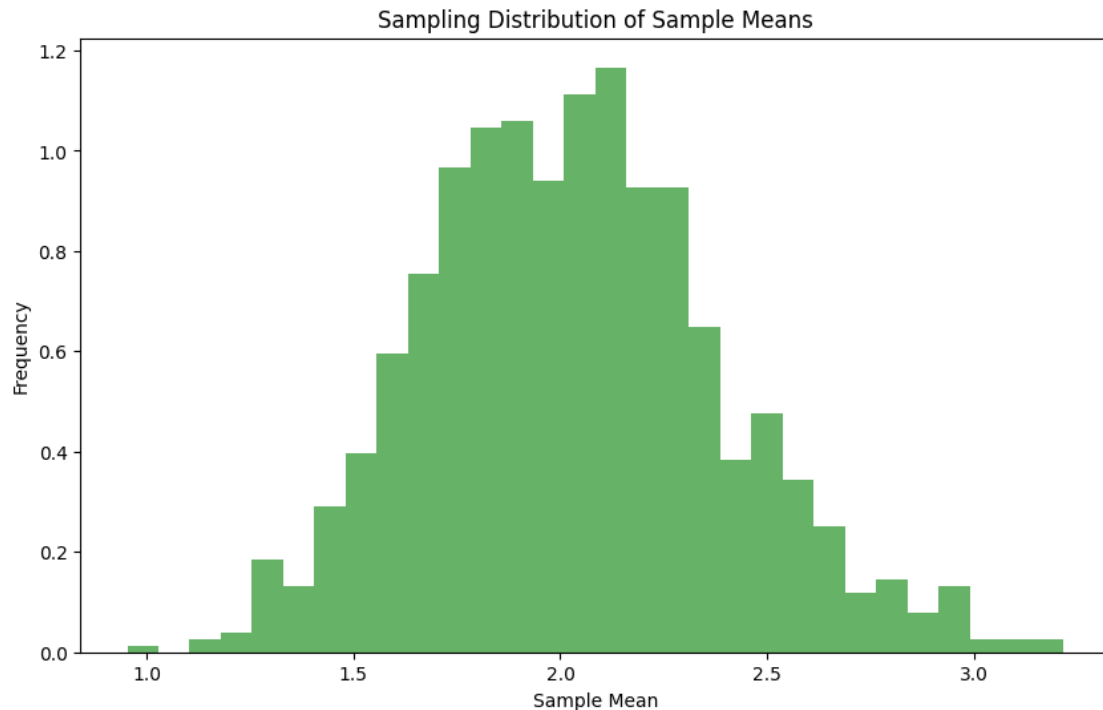
Graphs and visualizations play a crucial role in understanding and presenting data. Here are some example visualizations for the concepts discussed:

12 1. Sampling Distribution (Central Limit Theorem)

```
[9]: # Simulate population data
population = np.random.exponential(scale=2, size=1000)

# Draw samples and compute means
sample_means = [np.mean(np.random.choice(population, 30)) for _ in range(1000)]

# Plotting
plt.figure(figsize=(10, 6))
plt.hist(sample_means, bins=30, density=True, alpha=0.6, color='g')
plt.title('Sampling Distribution of Sample Means')
plt.xlabel('Sample Mean')
plt.ylabel('Frequency')
plt.show()
```



13 2. Confidence Interval Visualization

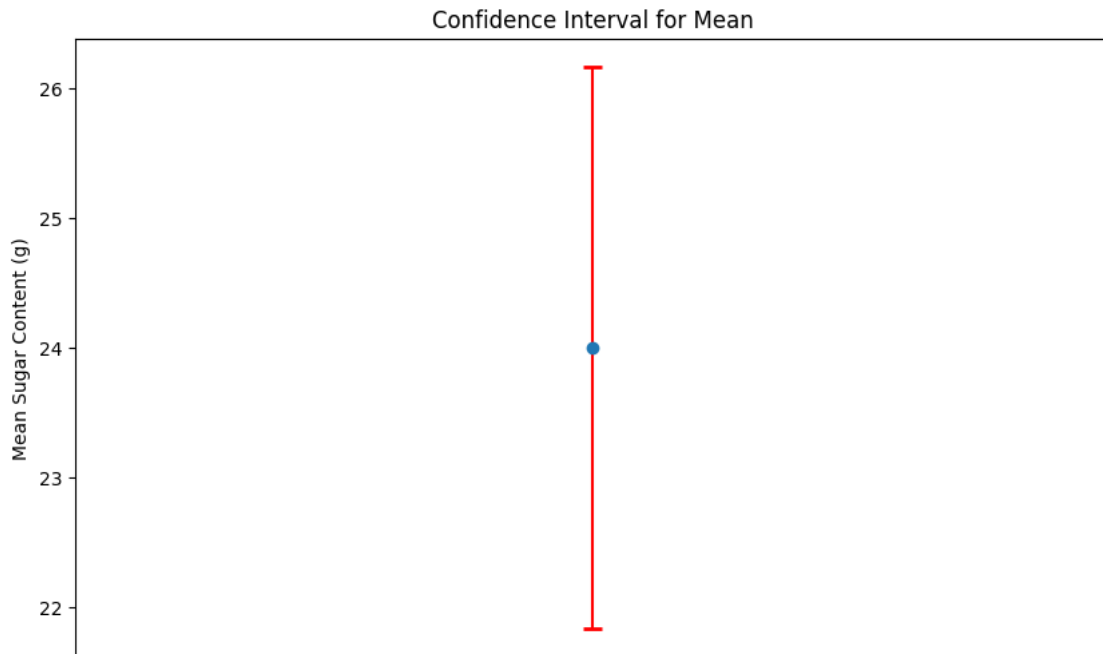
```
[10]: # Data for confidence interval
sample_mean = 24
std_dev = 8
sample_size = 37
confidence_level = 0.90

# Calculate margin of error and confidence interval
z_score = stats.norm.isf((1 - confidence_level) / 2)
margin_of_error = z_score * (std_dev / np.sqrt(sample_size))
confidence_interval = (sample_mean - margin_of_error, sample_mean +
    ↪margin_of_error)

# Plotting
plt.figure(figsize=(10, 6))
plt.errorbar(x=0, y=sample_mean, yerr=margin_of_error, fmt='o', capsize=5,
    ↪capthick=2, ecolor='red')
plt.title('Confidence Interval for Mean')
plt.ylabel('Mean Sugar Content (g)')
plt.xticks([])
```



```
plt.show()
```



Below are examples of bell curve graphs (normal distribution curves) illustrating the rejection regions for the null hypothesis H_0 in hypothesis testing. The rejection regions are based on the significance level α .

I'll provide visual representations for:

Left-Tailed Test Right-Tailed Test Two-Tailed Test Let's generate these graphs using Python.

Left-Tailed Test In a left-tailed test, the rejection region is on the left side of the mean. We reject H_0 if the test statistic falls in this region.

Right-Tailed Test In a right-tailed test, the rejection region is on the right side of the mean. We reject H_0 if the test statistic falls in this region.

Two-Tailed Test In a two-tailed test, there are two rejection regions: one on the left side and one on the right side of the mean. We reject H_0 if the test statistic falls in either of these regions.

Here are the bell curve graphs illustrating the rejection regions for hypothesis testing:

Left-Tailed Test:

The rejection region is on the left side of the mean. If the test statistic falls in this region, we reject the null hypothesis H_0 . Right-Tailed Test:

The rejection region is on the right side of the mean. If the test statistic falls in this region, we reject the null hypothesis H_0 . Two-Tailed Test:

There are two rejection regions, one on the left and one on the right side of the mean. If the test statistic falls in either of these regions, we reject the null hypothesis H_0 . These graphs visually demonstrate how the significance level α determines the critical values and rejection regions for different types of hypothesis tests.

14 Summary of Inferential Statistics in Python

Inferential statistics involves making predictions or inferences about a population based on a sample of data drawn from that population. This involves several key concepts and methods:

Parameters and Statistics:

Parameter: A measure for the entire population (e.g., population mean μ). Statistic: A measure from a sample of the population (e.g., sample mean \bar{x}). Sampling Techniques:

Simple Random Sampling: Every member has an equal chance of being selected. Stratified Sampling: The population is divided into strata, and samples are taken from each. Systematic Sampling: Selecting every k -th individual from the population. Central Limit Theorem (CLT):

States that the distribution of sample means approximates a normal distribution as the sample size becomes large. Estimation:

Point Estimation: Provides a single value estimate of a population parameter. Interval Estimation: Provides a range within which the population parameter is expected to lie, usually expressed with a confidence interval. Confidence Intervals:

Used to estimate the population parameter. Confidence intervals are calculated using Z-distribution or T-distribution depending on whether the population standard deviation is known and the sample size. Hypothesis Testing:

Null Hypothesis (H_0): The statement being tested, typically posits no effect or no difference. Alternative Hypothesis (H_a): The statement we want to test against the null hypothesis. Hypothesis Testing Methods:

Critical Value Approach: Compare the test statistic to a critical value. P-Value Approach: Calculate the p-value and compare it to the significance level α . Confidence Interval Approach: Check if the confidence interval contains the null hypothesis value. Rejection Regions:

Left-Tailed Test: The rejection region is in the left tail of the distribution. Right-Tailed Test: The rejection region is in the right tail of the distribution. Two-Tailed Test: There are rejection regions in both tails of the distribution. Visuals The provided graphs show the rejection regions for different types of tests based on a significance level ($\alpha = 0.05$):

Left-Tailed Test:

The rejection region is on the left side of the mean. Critical value: -1.64 Left-Tailed Test:

The rejection region is on the right side of the mean. Critical value: 1.64 Right-Tailed Test:

Rejection regions on both sides of the mean. Critical values: -1.96 and 1.96 Two-Tailed Test: These concepts and methods form the foundation of inferential statistics and are essential for analyzing data and making informed decisions based on sample data.

15 *Insights Gained from Inferential Statistics*

Understanding Population from Sample:

By studying a sample, we can infer characteristics about the entire population, which is often impractical or impossible to measure directly. Sampling Techniques:

Different sampling methods (simple random, stratified, systematic) help ensure that the sample represents the population, reducing bias and improving the reliability of inferences. Central Limit Theorem (CLT):

The CLT justifies the use of the normal distribution in many statistical procedures, even when the underlying data is not normally distributed, provided the sample size is sufficiently large. Estimations:

Point Estimation: Provides a single value estimate of a population parameter. Interval Estimation: More informative as it gives a range within which the population parameter is expected to lie, with a certain level of confidence (e.g., 95%). Confidence Intervals:

Confidence intervals offer a way to estimate population parameters with an associated level of certainty, helping to quantify the uncertainty inherent in sample data. Hypothesis Testing:

A structured approach to testing assumptions (hypotheses) about population parameters. It helps in making decisions or inferences about the population based on sample data. Types of Tests:

Left-Tailed Test: Tests if the parameter is less than a certain value. Right-Tailed Test: Tests if the parameter is greater than a certain value. Two-Tailed Test: Tests if the parameter is significantly different from a certain value (either less or more). Rejection Regions:

Visualizing rejection regions helps understand where the test statistic must lie to reject the null hypothesis. Rejection regions depend on the chosen significance level (α) and the type of test being performed. Critical Values and P-Values:

Critical values help determine the cutoff points for rejecting the null hypothesis. P-values provide the probability of obtaining test results at least as extreme as the results actually observed, under the assumption that the null hypothesis is correct. Error Types:

Type I Error: Rejecting the null hypothesis when it is true (false positive). Type II Error: Failing to reject the null hypothesis when it is false (false negative). Applications of Insights Business: Making data-driven decisions, like estimating the average customer spending or testing the effectiveness of a new marketing campaign. Healthcare: Estimating average recovery times and testing new treatments or drugs. Manufacturing: Quality control by estimating defect rates and testing process improvements. Social Sciences: Studying population behaviors, attitudes, and characteristics through sample surveys. Conclusion Inferential statistics is a powerful tool for making

predictions and informed decisions based on sample data. Understanding and applying these concepts can lead to more accurate and reliable conclusions in various fields, enhancing our ability to make evidence-based decisions.