Statistical Modeling - 24DS636 (2024-25)

Abhijith M S

2025-02-19

Table of contents

| Co | ourse | Introdu | | J |
|----|--------|----------|---|----|
| | | Syllab | us | 1 |
| | | Evalua | ations: A Tentative Timeline | 1 |
| 1 | Desc | criptive | statistics | 3 |
| | | 1.0.1 | Describing Data sets | 3 |
| | | 1.0.2 | Frequency Tables and Graphs | 3 |
| | | 1.0.3 | Relative Frequency Tables and Graphs | Ę |
| | | 1.0.4 | Grouped Data, Histograms, Ogives, and Stem and Leaf Plots | 6 |
| 2 | Para | ameter l | Estimation | 11 |
| | 2.1 | Point | Estimation: Maximum Likelihood Estimators | 11 |
| | 2.2 | Interv | al Estimates | 11 |
| | | 2.2.1 | Confidence Intervals for the Mean of a normal population with | |
| | | | known Variance | 11 |
| | | 2.2.2 | Confidence Intervals for the Mean of a normal population with | |
| | | | unknown Variance | 17 |
| | | 2.2.3 | Confidence Intervals for the Variance of a Normal Distribution $$ | 18 |
| 3 | Нур | othesis | Testing | 21 |
| | | Introd | uction | 23 |
| | | Testin | g a Null Hypothesis | 24 |
| | | Signifi | cance Level and Classical Approach | 25 |
| 4 | Нур | othesis | Tests Concerning the mean of a normal population | 27 |
| | 4.1 | With | known Variance (Z-test) | 27 |
| | | 4.1.1 | Choosing the Significance Level | 29 |
| | | 4.1.2 | Hypothesis Testing Summary: Z- Test | 29 |
| | 4.2 | With | unknown Variance (T-test) | 31 |
| | | 4.2.1 | Hypothesis Testing Summary: T- Test | 32 |
| 5 | Нур | othesis | Tests Concerning the variance of a normal population | 35 |
| | | 5.0.1 | Hypothesis Testing Summary: chi-square Test | 36 |
| Re | eferen | ces | | 39 |

Course Introduction

The website contains course contents of "Statistical Modeling" offered by Abhijith M S, PhD to Masters students pursuing M.Tech in Data Science, during the even semester of the academic year 2024-25.

Syllabus

(As given in the curriculum)

- Probability, Random Variables & Probability Distributions.
- Sampling, analysis of sample data-Empirical Distributions, Sampling from a Population Estimation, confidence intervals, point estimation—Maximum Likelihood, Probability mass functions, Modeling distributions, Hypothesis testing- Z, t, Chi-Square.
- ANOVA & Designs of Experiments Single, Two factor ANOVA, Factorials ANOVA models.
- Linear least squares, Correlation & Regression Models-linear regression methods, Ridge regression, LASSO, univariate and Multivariate Linear Regression, probabilistic interpretation, Regularization, Logistic regression, locally weighted regression
- $\bullet~$ Exploratory data analysis, Time series analysis, Analytical methods ARIMA and SARIMA.

Evaluations: A Tentative Timeline

- Best two marks out of three quizzes (Total = 20 marks)
- Quiz-1 (10 marks): (January First week)
- Quiz-2 (10 marks):(March First week)
- Quiz-3 (10 marks):(April First week)
- Assignments (Total = 30 marks)
- Assignment-1 (10 marks):(Submission: End of January)
- Assignment-2 (10 marks):(Submission: End of March)

Course Introduction

- Project Review 1 (10 marks):(February second week)
- Mid Sem (Total = 20 marks)
- Mid-Semester Exam (20 marks):(Feb first week, as per Academic calender)
- End Sem (Total = 30 marks)

 $Contact: \ ms_abhijith@cb.amrita.edu$

1 Descriptive statistics

- Descriptive statistics deals with methods to describe and summarize data.
- Describing of data is effectively done through tables or graphs. Those often reveal important features such as the range, the degree of concentration, and the symmetry of the data.
- The summary of data is expressed through numerical quantities (summary statistics) whose values are determined by the data.

1.0.1 Describing Data sets

1.0.2 Frequency Tables and Graphs

- A data set having a relatively small number of distinct values can be conveniently presented in a frequency table.
- Data from a frequency table can be graphically represented by:
 - Line Graph
 - Bar Graph
 - Frequency Polygon

```
import numpy as np
import matplotlib.pyplot as plt

# Sample data
data = np.array(['A', 'B', 'C', 'A', 'A', 'B','B','B','B','B','C','C','C','C'])

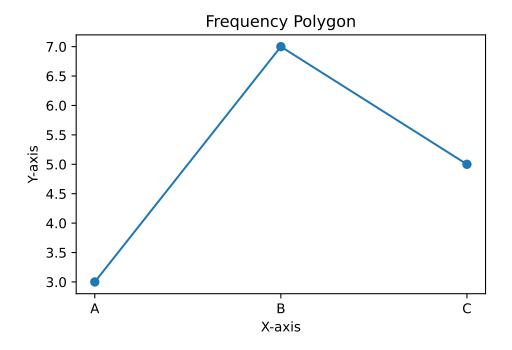
# Calculate frequencies
values, frequencies = np.unique(data, return_counts=True)

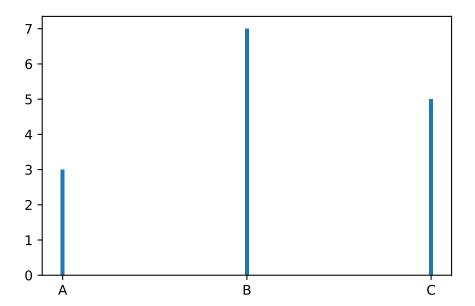
# Line Graph
plt.plot(values, frequencies, marker='o')
plt.title('Frequency Polygon')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
```

1 Descriptive statistics

```
plt.show()

plt.bar(values, frequencies, width=0.02)
plt.show()
```





1.0.3 Relative Frequency Tables and Graphs

- Consider a data set consisting of n values. If f is the frequency of a particular value, then the ratio f /n is called its relative frequency.
- That is, the relative frequency of a data value is the proportion of the data that have that value.
- The relative frequencies can be represented graphically by:
 - relative frequency line
 - relative frequency bar graph
 - relative frequency polygon
 - pie chart: A pie chart is often used to indicate relative frequencies when the
 data are not numerical in nature. A circle is constructed and then sliced into
 different sectors with areas proportional to the respective relative frequencies.

```
import numpy as np
import matplotlib.pyplot as plt

# Sample data
data = np.array(['A', 'B', 'C', 'A', 'A', 'B','B','B','B','B','B','C','C','C','C'])

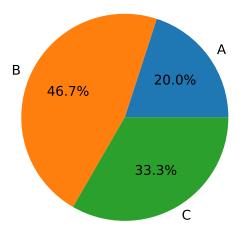
# Calculate frequencies
values, frequencies = np.unique(data, return_counts=True)

relative_frequencies = frequencies/len(data)
print(relative_frequencies)

plt.pie(relative_frequencies, labels = values, autopct='%1.1f%%')
plt.show()
```

 $[0.2 \quad 0.46666667 \ 0.333333333]$

1 Descriptive statistics



1.0.4 Grouped Data, Histograms, Ogives, and Stem and Leaf Plots

- For some data sets the number of distinct values is too large to utilize frequency tables.
- Instead, in such cases, it is useful to divide the values into groupings, or class intervals, and then plot the number of data values falling in each class interval.
- The number of class intervals chosen should be a trade-off between:
 - choosing too few classes at a cost of losing too much information about the actual data values in a class.
 - choosing too many classes, which will result in the frequencies of each class being too small.
- It is common, although not essential, to choose class intervals of equal length.
- The endpoints of a class interval are called the class boundaries.
- We will adopt the left-end inclusion convention, which stipulates that a class interval contains its left-end but not its right-end boundary point.
- Thus, for instance, the class interval 20-30 contains all values that are both greater than or equal to 20 and less than 30.
- A bar graph plot of class data, with the bars placed adjacent to each other, is called a histogram.
- The vertical axis of a histogram can represent either the class frequency or the relative class frequency.

```
import numpy as np
import matplotlib.pyplot as plt

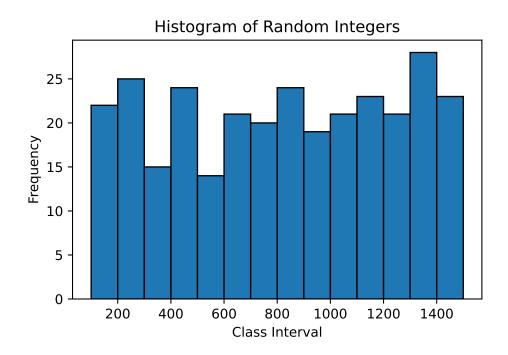
# Set seed for reproducibility
np.random.seed(50)

# Generate 300 random integers between 100 and 1500
random_integers = np.random.randint(100, 1501, size=300)

# Create bins for the class intervals
bins = np.arange(100, 1600, 100)
print(bins)

# Plot histogram
plt.hist(random_integers, bins=bins, edgecolor='black')
plt.title('Histogram of Random Integers')
plt.xlabel('Class Interval')
plt.ylabel('Frequency')
plt.show()
```

 $[\ 100\ \ 200\ \ 300\ \ 400\ \ 500\ \ 600\ \ 700\ \ 800\ \ 900\ \ 1000\ \ 1100\ \ 1200\ \ 1300\ \ 1400$ 1500]

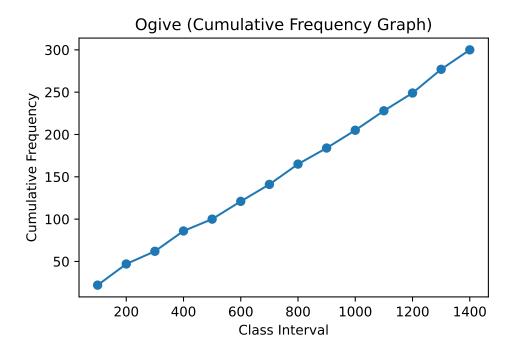


1 Descriptive statistics

- We are sometimes interested in plotting a cumulative frequency (or cumulative relative frequency) graph.
- A point on the horizontal axis of such a graph represents a possible data value; its corresponding vertical plot gives the number (or proportion) of the data whose values are less than or equal to it.
- A cumulative frequency plot is called an ogive.

```
import numpy as np
import matplotlib.pyplot as plt
# Set seed for reproducibility
np.random.seed(50)
# Generate 300 random integers between 100 and 1500
random integers = np.random.randint(100, 1501, size=300)
# Create bins for the class intervals
bins = np.arange(100, 1600, 100)
print(bins)
histograms = np.histogram(random_integers, bins=bins)[0]
print(histograms)
cumulativeSum = np.cumsum(histograms)
print(cumulativeSum)
plt.plot(bins[:-1], cumulativeSum, marker='o', linestyle='-')
plt.title('Ogive (Cumulative Frequency Graph)')
plt.xlabel('Class Interval')
plt.ylabel('Cumulative Frequency')
#plt.grid(True)
plt.show()
```

```
[ 100 200 300 400 500 600 700 800 900 1000 1100 1200 1300 1400 1500]
[22 25 15 24 14 21 20 24 19 21 23 21 28 23]
[ 22 47 62 86 100 121 141 165 184 205 228 249 277 300]
```



- An efficient way of organizing a small-to moderate-sized data set is to utilize a stem and leaf plot.
- Such a plot is obtained by first dividing each data value into two parts its stem and its leaf.
- For instance, if the data are all two-digit numbers, then we could let the stem part of a data value be its tens digit and let the leaf be its ones digit.
- Thus, for instance, the value 62 is expressed as

stem_leaf[stem].append(leaf)

else:

```
import numpy as np

# Sample data
data = np.array([62, 67, 63, 68, 69, 61, 64, 65, 66, 60, 75, 74, 76, 78, 90, 92, 34, 36, 56, 57, 45, 53, 52, 59, 73, 74

# Create stem and leaf plot
stem_leaf = {}

for number in data:
    stem = number // 10
    leaf = number % 10
    if stem in stem_leaf:
```

1 Descriptive statistics

```
stem_leaf[stem] = [leaf]

# Print stem and leaf plot
for stem, leaves in sorted(stem_leaf.items()):
    print(f"{stem} | {' '.join(map(str, sorted(leaves)))}")
```

2 Parameter Estimation

2.1 Point Estimation: Maximum Likelihood Estimators

(Please refer the book titled "Introduction to Probability and Statistics for Engineers and Scientists" by Sheldon M Ross for more details)

2.2 Interval Estimates

- Consider a sample X_1, X_2, \dots, X_n drawn from a known distribution with an unknown mean μ .
- It is established that the sample mean $\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$ serves as the maximum likelihood estimator for μ .
- However, the sample mean \bar{X} is not expected to be exactly equal to μ , but rather close to it.
- Therefore, instead of providing a single point estimate, it is often more useful to specify an interval within which we are confident that μ lies.
- To determine such an interval estimator, we utilize the probability distribution of the point estimator.

2.2.1 Confidence Intervals for the Mean of a normal population with known Variance

- Consider a sample $X_1, X_2, ..., X_n$ drawn from a normal distribution with an unknown mean μ and a known variance σ^2 .
- The point estimator \bar{X} is normal with mean μ and variance σ^2/n .
- Therefore, $\frac{\bar{X}-\mu}{\sigma/\sqrt{n}}$ follows a standard normal distribution.

$$P\left(-1.96 < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < 1.96\right) = 0.9750 - 0.0250 = 0.95$$

Parameter Estimation above equation:

i What to do

$$P\left(-1.96\frac{\sigma}{\sqrt{\pi}} < \bar{X} - \mu < 1.96\frac{\sigma}{\sqrt{\pi}}\right) = 0.95$$

What to do $P\left(-1.96\frac{\sigma}{\sqrt{n}} < \bar{X} - \mu < 1.96\frac{\sigma}{\sqrt{n}}\right) = 0.95$ Consider that I want to find an interval around \bar{X} such that the actual population mean μ falls within the interval, say 95% of the times. $P\left(1.96\frac{\pi}{\sqrt{n}} > \mu - X > -1.96\frac{\pi}{\sqrt{n}}\right) = 0.95$

$$P\left(-1.96\frac{\sigma}{\sqrt{n}} < \mu - \bar{X} < 1.96\frac{\sigma}{\sqrt{n}}\right) = 0.95$$

$$P\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

- We have P(Z < -1.96) = 0.025, similarly P(Z > 1.96) = 0.025. Usually 1.96 is represented generally as $z_{0.025}$. Thus, $P(Z < -z_{0.025}) = 0.025$ and P(Z > $z_{0.025}$) = 0.025.
- Hence, 100(1-0.05) percent confidence interval for the mean of a normal population with known variance is:

$$P\left(\bar{X} - z_{0.025} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{0.025} \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

$$P\left(\bar{X} - z_{0.05/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{0.05/2} \frac{\sigma}{\sqrt{n}}\right) = (1 - 0.05)$$

• For a confidence level of $100(1-\alpha)$ percent, the corresponding critical value from the standard normal distribution is $z_{\alpha/2}$.

• The $100(1-\alpha)$ percent confidence interval for μ is given by:

$$\mu \in \left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) \tag{2.1}$$

- The interval as given in Equation 2.1 is called a two-sided confidence interval.
- Also the term $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ is called the margin of error.

Derivation of two-sided confidence interval

• To find $100(1-\alpha)$ percent confidence interval of mean (μ) , we have;

$$P\left(-z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < z_{\alpha}/2\right) = 1 - \alpha$$

• Doing the same manipulations we did earlier for obtaining the 95percent confidence interval we can obtain:

$$P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

• The above equation give us the required confidence interval, as given in Equation 2.1.

What if!?

What if we are interested in one sided confidence intervals !!?

• One-sided Upper Confidence Iterval

• To determine such an interval, for a standard normal random variable Z, we have;

$$P(Z < 1.645) = 0.95$$

• Thus,

$$P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < 1.645\right) = 0.95$$

$$P\left(\mu - \bar{X} > -1.645 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

$$P\left(\mu > \bar{X} - 1.645 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

Thus a 95 percent one-sided upper confidence interval for μ is

$$\mu \in \left(\bar{X} - 1.645 \frac{\sigma}{\sqrt{n}}, \infty\right)$$

or in other words; 100(1-0.05) percent one-sided upper confidence interval for μ is

$$\mu \in \left(\bar{X} - z_{0.05} \frac{\sigma}{\sqrt{n}}, \infty\right)$$

Oneside interval!

Can you think of another one sided confidence interval?

- One-sided lower confidence interval
 - We have

$$P(Z > -1.645) = 0.95$$

Proceed just like in the previous case and you will find a 100(1-0.05) percent one-sided lower confidence interval for μ as;

$$\mu \in \left(-\infty, \bar{X} + z_{0.05} \frac{\sigma}{\sqrt{n}}\right)$$

• In general, $100(1-\alpha)$ percent one-sided upper confidence interval for μ is given in Equation 2.2.

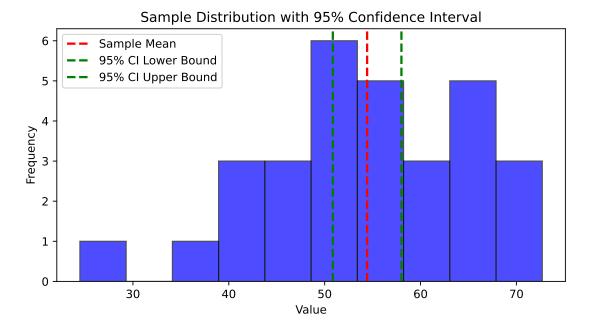
$$\mu \in \left(\bar{X} - z_{\alpha} \frac{\sigma}{\sqrt{n}}, \infty\right)$$
 (2.2)

• Also, $100(1-\alpha)$ percent one-sided lower confidence interval for μ is given in Equation 2.3.

$$\mu \in \left(-\infty, \bar{X} + z_{\alpha} \frac{\sigma}{\sqrt{n}}\right) \tag{2.3}$$

• The python code below creates a sample and find 95% confidence interval for the mean if the population standard deviation is assumed to be 10. Other values are specified in the code.

```
import numpy as np
import matplotlib.pyplot as plt
# Parameters
mu = 50 \# true mean
sigma = 10 \# known standard deviation
n = 30 \# sample size
alpha = 0.05 \# significance level
# Generate a sample
np.random.seed(0)
sample = np.random.normal(mu, sigma, n)
sample\_mean = np.mean(sample)
# Calculate the confidence interval
z = 1.96 \# z-value for 95% confidence
margin\_of\_error = z * (sigma / np.sqrt(n))
confidence_interval = (sample_mean - margin_of_error, sample_mean + margin_of_error)
# Plot the sample and confidence interval
plt.figure(figsize=(8, 4))
plt.hist(sample, bins=10, alpha=0.7, color='blue', edgecolor='black')
plt.axvline(sample_mean, color='red', linestyle='dashed', linewidth=2, label='Sample Mean')
plt.axvline(confidence_interval[0], color='green', linestyle='dashed', linewidth=2, label='95% CI Lower Bound')
plt.axvline(confidence_interval[1], color='green', linestyle='dashed', linewidth=2, label='95% CI Upper Bound'
plt.title('Sample Distribution with 95% Confidence Interval')
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.legend()
plt.show()
print(f"Sample Mean: {sample_mean}")
print(f"95% Confidence Interval: {confidence_interval}")
```



Sample Mean: 54.42856447263174

95% Confidence Interval: (50.85011043026466, 58.007018514998826)

Problem

Suppose that when a signal having value μ is transmitted from location A the value received at location B is normally distributed with mean μ and variance 4. That is, if μ is sent, then the value received is $\mu + N$ where N, representing noise, is normal with mean 0 and variance 4. To reduce error, suppose the same value is sent 9 times. If the successive values received are 5, 8.5, 12, 15, 7, 9, 7.5, 6.5, 10.5; (a). construct a 95 percent two-sided confidence interval for μ .

(b). construct 95 percent one-sided upper and lower confidence intervals for μ .

Problem

Suppose a quality control manager at a factory wants to ensure that the average weight of a product is at least 500 grams. They take a random sample of 30 products and find the sample mean weight to be 495 grams with a standard deviation of 10 grams. Help the manager to estimate the minimum average weight of the products with 95% confidence.

- 2.2.2 Confidence Intervals for the Mean of a normal population with unknown Variance
 - If you recollect the discussion we had about the sample mean from a normal population with unknown variance we saw that variable t_{n-1} given by:

$$t_{n-1} = \sqrt{n} \frac{\bar{X} - \mu}{S}$$

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

• Because of the symmetry of the t-distribution we can write for any $\alpha \in (0, 1/2)$;

$$\begin{split} P\left(-t_{\alpha/2,n-1} < \sqrt{n}\frac{\bar{X}-\mu}{S} < t_{\alpha/2,n-1}\right) &= 1-\alpha \\ \\ P\left(-\bar{X}-t_{\alpha/2,n-1}\frac{\sqrt{n}}{S} < -\mu < -\bar{X}+t_{\alpha/2,n-1}\frac{\sqrt{n}}{S}\right) &= 1-\alpha \\ \\ P\left(\bar{X}+t_{\alpha/2,n-1}\frac{\sqrt{n}}{S} > \mu > \bar{X}-t_{\alpha/2,n-1}\frac{\sqrt{n}}{S}\right) &= 1-\alpha \\ \\ P\left(\bar{X}-t_{\alpha/2,n-1}\frac{S}{\sqrt{n}} < \mu < \bar{X}+t_{\alpha/2,n-1}\frac{S}{\sqrt{n}}\right) &= 1-\alpha \end{split}$$

• If the sample mean is \bar{X} and sample standard deviation S, then we can say that with $100(1-\alpha)$ percent confidence that

$$\mu \in \left(\bar{X} - t_{\alpha/2,n-1}\frac{S}{\sqrt{n}}, \bar{X} + t_{\alpha/2,n-1}\frac{S}{\sqrt{n}}\right)$$

• In this case $100(1-\alpha)$ percent one-sided upper confidence interval can be obtained from the fact that:

$$P\left(\sqrt{n}\frac{(\bar{X}-\mu)}{S} < t_{\alpha,n-1}\right) = 1 - \alpha$$

$$P\left(\mu > \bar{X} - \frac{S}{\sqrt{n}}t_{\alpha,n-1}\right) = 1 - \alpha$$

• Thus $100(1 - \alpha)$ percent one-sided upper confidence interval for the mean in this case is given by;

$$\mu \in \left(\bar{X} - \frac{S}{\sqrt{n}} t_{\alpha,n-1}, \infty\right)$$

• Thus $100(1-\alpha)$ percent one-sided lower confidence interval for the mean in this case is given by;

$$\mu \in \left(-\infty, \bar{X} + \frac{S}{\sqrt{n}}t_{\alpha, n-1}\right)$$

🍐 Problem

Let us again consider the previous problem but let us now suppose that when the value μ is transmitted at location A then the value received at location B is normal with mean μ and variance σ^2 but with σ^2 being unknown. If 9 successive values are, 5, 8.5, 12, 15, 7, 9, 7.5, 6.5, and 10.5, compute a 95 percent confidence interval for μ .

2.2.3 Confidence Intervals for the Variance of a Normal Distribution

• If we are sampling from a normal distribution with unknown mean and unknown variance then;

$$(n-1)\frac{S^2}{\sigma^2} \sim \chi_{n-1}^2$$

follows a chi-squared distribution.

• We have

$$P\left(\chi^2_{1-\alpha/2,n-1} \le (n-1)\frac{S^2}{\sigma^2} \le \chi^2_{\alpha/2,n-1}\right) = 1 - \alpha$$

$$P\left(\chi^2_{1-\alpha/2,n-1} \leq (n-1)\frac{S^2}{\sigma^2} \leq \chi^2_{\alpha/2,n-1}\right) = 1-\alpha$$

$$P\left(\frac{(n-1)S^2}{\chi^2_{\alpha/2,n-1}} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi^2_{1-\alpha/2,n-1}}\right) = 1-\alpha$$

• Hence, $100(1-\alpha)$ percent two-sided confidence interval for the variance in this case;

$$\sigma^2 \in \left(\frac{(n-1)S^2}{\chi^2_{\alpha/2,n-1}}, \frac{(n-1)S^2}{\chi^2_{1-\alpha/2,n-1}}\right)$$

• The $100(1-\alpha)$ percent one-sided upper and lower confidence intervals in this case will be respectively;

$$\left(\frac{(n-1)S^2}{\chi^2_{\alpha,n-1}},\infty\right)$$

and

$$\left(0, \frac{(n-1)S^2}{\chi^2_{1-\alpha, n-1}}\right)$$

Problem

A standardized procedure is expected to produce washers with very small deviation in their thicknesses. Suppose that 10 such washers were chosen and measured. If the thicknesses of these washers were, in inches; .123, .133, .124, .125, .126, .128, .120, .124, .130, and .126. What is a 90 percent confidence interval for the standard deviation of the thickness of a washer produced by this procedure?

All problems and most part of text are taken from Ross (2009).

Hypothesis Testing

Problem

A quality control manager at a factory wants to ensure that the average weight of a product is at least 500 grams. They take a random sample of 30 products and find the sample mean weight to be 495 grams with a standard deviation of 10 grams. The manager wants to test if the average weight of the products is significantly less than 500 grams at a 5% significance level.

- Null Hypothesis (H_0) : The average weight of the products is at least 500 grams ($\mu \geq 500$).
- Alternative Hypothesis (H_1) : The average weight of the products is less than 500 grams ($\mu < 500$).

Examples Highlighting the Need for Hypothesis Testing

1. Medical Research:

- Scenario: A pharmaceutical company develops a new drug intended to lower blood pressure.
- Hypothesis Testing: The null hypothesis (H_0) might state that the new drug has no effect on blood pressure, while the alternative hypothesis (H_1) states that the drug does lower blood pressure. Hypothesis testing helps determine if the observed effects in clinical trials are statistically significant or if they could have occurred by random chance.

2. Quality Control:

- Scenario: A factory produces light bulbs, and the quality control team wants to ensure that the average lifespan of the bulbs is 1000 hours.
- Hypothesis Testing: The null hypothesis (H_0) could be that the mean lifespan of the bulbs is 1000 hours. The alternative hypothesis (H_1) might be that the mean lifespan is not 1000 hours. Hypothesis testing helps the team decide whether to accept the production process or take corrective actions.

3. Marketing:

- Scenario: A company launches a new advertising campaign and wants to know if it has increased sales.
- Hypothesis Testing: The null hypothesis (H_0) might state that the advertising campaign has no effect on sales, while the alternative hypothesis (H_1) states that the campaign has increased sales. Hypothesis testing helps the company determine if the increase in sales is statistically significant.

4. Education:

- Scenario: An educator wants to test if a new teaching method is more effective than the traditional method.
- Hypothesis Testing: The null hypothesis (H_0) could be that there is no difference in effectiveness between the new and traditional methods. The alternative hypothesis (H_1) might be that the new method is more effective. Hypothesis testing helps in making data-driven decisions about adopting new teaching strategies.

5. Environmental Science:

- Scenario: Researchers want to determine if a new policy has reduced pollution levels in a city.
- Hypothesis Testing: The null hypothesis (H_0) might state that the policy has no effect on pollution levels, while the alternative hypothesis (H_1) states that the policy has reduced pollution levels. Hypothesis testing helps in evaluating the effectiveness of environmental policies.

These examples illustrate how hypothesis testing is a crucial tool in various fields for making informed decisions based on data.

Important Terminology

- Null Hypothesis (H_0) : The hypothesis that there is no effect or no difference. It is the default assumption that any observed effect is due to random chance. It is the hypothesis that researchers aim to test against.
- Alternative Hypothesis (H_1 or H_a): The hypothesis that there is an effect or a difference. It is what researchers want to prove.
- Test Statistic: A standardized value that is calculated from sample data during a hypothesis test. It is used to decide whether to reject the null hypothesis.
- P-value: The probability of obtaining test results at least as extreme as the observed results, assuming that the null hypothesis is true. A smaller p-value

indicates stronger evidence against the null hypothesis.

- Significance Level (α): A threshold set by the researcher which the p-value must be below in order to reject the null hypothesis. Common significance levels are 0.05, 0.01, and 0.10.
- Critical Value: The value that the test statistic must exceed in order to reject the null hypothesis. It is determined based on the significance level and the distribution of the test statistic.
- Power of a Test: The probability that the test correctly rejects a false null hypothesis (i.e., it does not make a type II error). Higher power indicates a greater ability to detect an effect when there is one.
- Type I Error: The error made when the null hypothesis is true, but is incorrectly rejected. The probability of making a type I error is denoted by α .
- Type II Error: The error made when the null hypothesis is false, but is incorrectly accepted. The probability of making a type II error is denoted by β .
- Confidence Interval: A range of values derived from the sample data that is likely to contain the population parameter. It provides an estimate of the parameter with a certain level of confidence (e.g., 95%).
- One-tailed Test: A hypothesis test in which the region of rejection is on only one side of the sampling distribution. It tests for the possibility of the relationship in one direction.
- Two-tailed Test: A hypothesis test in which the region of rejection is on both sides of the sampling distribution. It tests for the possibility of the relationship in both directions.

Introduction

- A statistical hypothesis is typically a statement regarding a set of parameters of a population distribution.
- It is termed a hypothesis because its true value is unknown.
- The main challenge is to devise a method to determine whether the values of a random sample from this population align with the hypothesis.
- Consider a population with distribution F_{θ} , where θ is unknown.
- We aim to test a specific hypothesis about θ .

3 Hypothesis Testing

- This hypothesis is denoted by H_0 and is referred to as the null hypothesis.
- For instance, if F_{θ} is a normal distribution function with mean θ and variance equal to 1, two possible null hypotheses about θ are:

$$H_0: \theta = 1$$

$$H_0: \theta > 1$$

$$H_0: \theta \leq 1$$

- It is important to note that the null hypothesis in the first case fully specifies the population distribution.
- Whereas the null hypothesis in the second and third cases do not.

i Simple and Composite Hypotheses

- A hypothesis that fully specifies the population distribution when true is known as a simple hypothesis. (Eg; $H_0:\theta=1$)
- A hypothesis that does not fully specifies the population distribution is referred to as a composite hypothesis. (Eg; $H_0: \theta > 1$, H_{0}: 1)

Testing a Null Hypothesis

- To test a specific null hypothesis H_0 , we take a sample of size n from the population, say X_1, X_2, \dots, X_n .
- Based on these n values, we decide whether to accept or reject H_0 .
- We define a region C in the n-dimensional space. This region is called the critical region.
- If the sample X_1, X_2, \dots, X_n falls within the critical region C, we reject H_0 . Otherwise, we accept H_0 .
- In simple terms, the critical region C helps us determine the outcome of the statistical test.

accepts
$$H_0$$
 if $(X_1, X_2, ..., X_n) \notin C$

and

rejects
$$H_0$$
 if $(X_1, X_2, ..., X_n) \in C$

Types of Errors in Hypothesis Testing

- When developing a procedure for testing a given null hypothesis H_0 , it is crucial to recognize that two different types of errors can occur.
- A type I error occurs if the test incorrectly rejects H_0 when it is actually true.
- A type II error occurs if the test incorrectly accepts H_0 when it is actually false

Note

The goal of a statistical test for H_0 is not to definitively determine its truth but to assess if the data is consistent with H_0 .

Significance Level and Classical Approach

- H_0 should be rejected only if the observed data is highly unlikely under H_0 .
- The classical method involves specifying a value α , known as the level of significance.
- The test is designed so that the probability of rejecting H_0 when it is true does not exceed α .
- Common choices for α are 0.1, 0.05, and 0.005.
- This approach ensures that the probability of a type I error (incorrectly rejecting H_0) is controlled and does not exceed the chosen α .

Example

- For instance, consider testing the hypothesis that the mean of a normal distribution with parameters $(\theta, 1)$ is equal to 1.
- The test rejects the null hypothesis if the point estimate of θ (i.e., the sample mean) deviates more than $\frac{1.96}{\sqrt{n}}$ from 1.

3 Hypothesis Testing

• As we will discuss in the next section, the value $\frac{1.96}{\sqrt{n}}$ is selected to achieve a significance level of $\alpha=0.05$.

4 Hypothesis Tests Concerning the mean of a normal population

4.1 With known Variance (Z-test)

- Let $X_1, X_2, ..., X_n$ be a sample of size n from a normal distribution with an unknown mean μ and a known variance σ^2 .
- We are interested in testing the null hypothesis:

$$H_0: \mu = \mu_0$$

• Against the alternative hypothesis:

$$H_1: \mu \neq \mu_0$$

- Where μ_0 is a specified constant.
- Since $\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$ is a natural point estimator of μ , it is reasonable to accept H_0 if \bar{X} is not too far from μ_0 .
- Thus, the critical region of the test would be of the form:

$$C=\{X_1,\dots,X_n:|\bar{X}-\mu_0|>c\}$$

for some suitably chosen value c.

• To ensure that the test has a significance level α , we must determine the critical value c in the above equation such that the type I error is equal to α . This means c must satisfy:

$$P_{\mu_0}\{|\bar{X}-\mu_0|>c\}=\alpha$$

where P_{μ_0} denotes that the probability is computed under the assumption that population mean, $\mu = \mu_0$.

- 4 Hypothesis Tests Concerning the mean of a normal population
 - When $\mu = \mu_0$, \bar{X} follows a normal distribution with mean μ_0 and variance $\frac{\sigma^2}{n}$. Therefore, the standardized variable Z defined by:

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

will have a standard normal distribution.

• The probability of a type I error is given by:

$$P\left(|\bar{X} - \mu_0| > c\right) = \alpha$$

• Equivalently, this can be written as:

$$2P\left(Z > \frac{c\sqrt{n}}{\sigma}\right) = \alpha$$

• Where Z is a standard normal random variable. We know that:

$$P\left(Z > z_{\alpha/2}\right) = \frac{\alpha}{2}$$

• Therefore, we have:

$$\frac{c\sqrt{n}}{\sigma} = z_{\alpha/2}$$

• Solving for c, we get:

$$c = \frac{z_{\alpha/2}\sigma}{\sqrt{n}}$$

• Thus, the test at significance level α is to reject H_0 if:

$$|\bar{X} - \mu_0| > \frac{z_{\alpha/2}\sigma}{\sqrt{n}}$$

- And accept H_0 otherwise. Equivalently, we can reject H_0 if:

$$\sqrt{n}\frac{|\bar{X}-\mu_0|}{\sigma}>z_{\alpha/2}$$

• And accept H_0 if:

$$\sqrt{n} \frac{|\bar{X} - \mu_0|}{\sigma} \le z_{\alpha/2}$$

🌢 Problem

If a signal of value μ is sent from location A, then the value received at location B is normally distributed with mean μ and standard deviation 2. That is, the random noise added to the signal is an N(0, 4) random variable. There is reason for the people at location B to suspect that the signal value $\mu = 8$ will be sent today. Test this hypothesis if the same signal value is independently sent five times and the average value received at location B is X = 9. 5.

4.1.1 Choosing the Significance Level

- The appropriate significance level α depends on the specific context and consequences of the hypothesis test.
- If rejecting the null hypothesis H_0 would lead to significant costs or consequences, a more conservative significance level (e.g., 0.05 or 0.01) should be chosen.
- If there is a strong initial belief that H_0 is true, strict evidence is required to reject H_0 , implying a lower significance level.
- The test can be described as follows: For an observed value of the test statistic $\sqrt{n} \frac{|\bar{X} - \mu_0|}{\sigma}$, denoted as v, reject H_0 if the probability of the test statistic being as large as v under H_0 is less than or equal to α .
- This probability is known as the p-value of the test. H_0 is accepted if α is less than the p-value and rejected if α is greater than or equal to the p-value.
- In practice, the significance level is sometimes not set in advance. Instead, the p-value is calculated from the data, and decisions are made based on the p-value.
- If the p-value is much larger than any reasonable significance level, H_0 is accepted. Conversely, if the p-value is very small, H_0 is rejected.

4.1.2 Hypothesis Testing Summary: Z- Test

Table 4.1: Caption: Summary of hypothesis testing for a sample from a $N(\mu, \sigma^2)$ population with known σ^2 .

| Sample and Population | Details |
|-----------------------|----------------------|
| Sample | $\{X_1, X_2,, X_n\}$ |

4 Hypothesis Tests Concerning the mean of a normal population

| Sample and Population | Details |
|---|---|
| Population Known Parameter Sample Mean Significance Level | $N(\mu, \sigma^2)$ σ^2 $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ α |

| Hypothesis | Test Statistic (TS) | Reject if | p-Value if $TS = t$ |
|---|----------------------------------|-----------------------|---------------------|
| $H_0: \mu = \mu_0 \text{ vs}$ | $\sqrt{n}(\bar{X}-\mu_0)/\sigma$ | $ TS > z_{\alpha/2}$ | $2P\{Z \ge t \}$ |
| $H_1: \mu \neq \mu_0$ | | T. C | D(G) |
| 0 0 | $\sqrt{n}(\bar{X}-\mu_0)/\sigma$ | $TS > z_{\alpha}$ | $P\{Z \ge t\}$ |
| $H_1: \mu > \mu_0$ $H_0: \mu > \mu_0 \text{ vs}$ | $\sqrt{n}(\bar{X}-\mu_0)/\sigma$ | $TS < -z_{\alpha}$ | $P\{Z \le t\}$ |
| $H_1: \mu < \mu_0$ | • () ()// | a | |

♦ Problem

Imagine you're the quality control manager at a company that prides itself on the precision of its product weights. The company claims that the average weight of their product is exactly 100 grams. But, as a diligent manager, you decide to put this claim to the test. You randomly select a sample of 30 products and measure their weights. To your surprise, the average weight of your sample is 110 grams! Now, you need to determine if this difference is statistically significant or just a fluke. Assume the population standard deviation as 15 grams.

```
import numpy as np
from scipy import stats

# Given data
sample_mean = 495
population_mean = 500
std_dev = 10
sample_size = 30
alpha = 0.05

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

# Calculate the p-value
p_value = stats.norm.cdf(z_score)
```

```
# Determine if we reject the null hypothesis
reject_null = p_value < alpha

# Output the results
print(f"Z-score: {z_score}")
print(f"P-value: {p_value}")
print(f"Reject the null hypothesis: {reject_null}")</pre>
```

Z-score: -2.7386127875258306 P-value: 0.00308494966027208 Reject the null hypothesis: True

4.2 With unknown Variance (T-test)

- Let $X_1, X_2, ..., X_n$ be a sample of size n from a normal distribution with an unknown mean μ and a unknown variance.
- Say, we are interested in testing the null hypothesis:

$$H_0: \mu = \mu_0$$

• Against the alternative hypothesis:

$$H_1: \mu \neq \mu_0$$

- Where μ_0 is a specified constant.
- In the previous case (with known variance), for a significance level (α) we accepted the null hypothesis if:

$$\left| \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \right| \le z_{\alpha/2}$$

- But in this case, σ is unknown.
- We know that the statistic, T, as given below has a t-distribution with n-1 degrees of freedom when $\mu = \mu_0$.

$$T = \frac{\bar{X} - \mu_0}{S\sqrt{n}}$$

where S is the sample standard deviation.

- 4 Hypothesis Tests Concerning the mean of a normal population
 - Hence here with H_0 : $\mu=\mu_0$ and H_1 ; $\mu\neq\mu_0$; analogous to the Z-test here in T-test we can:

i two-sided t-test

• reject the null hypothesis (H_0) if:

$$\left|\frac{\bar{X}-\mu_0}{S/\sqrt{n}}\right| > t_{\alpha/2}$$

• accept H_0 if:

$$\left|\frac{\bar{X} - \mu_0}{S/\sqrt{n}}\right| \le t_{\alpha/2}$$

4.2.1 Hypothesis Testing Summary: T- Test

Table 4.3: Caption: Summary of hypothesis testing for a sample from a $N(\mu, \sigma^2)$ population with unknown σ^2 .

| P o P described | |
|-----------------------|--|
| Sample and Population | Details |
| Sample | $\{X_1, X_2,, X_n\}$ |
| Population | $N(\mu, \sigma^2)$ |
| Sample Mean | $\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$ |
| Sample Variance | $S^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (X_i - \bar{X})^2$ |
| Significance Level | α |

| Hypothesis | Test Statistic (TS) | Reject if | p-Value if $TS = t$ |
|--|--|-----------------------|-------------------------|
| $\overline{H_0: \mu = \mu_0 \text{ vs}}$ | $\sqrt{n}(\bar{X}-\mu_0)/S$ | TS > | $2P\{T_{n-1} \ge t \}$ |
| $H_1: \mu \neq \mu_0$ | | $t_{\alpha/2,n-1}$ | |
| | $\sqrt{n}(\bar{X}-\mu_0)/S$ | $TS > t_{\alpha,n-1}$ | $P\{T_{n-1} \ge t\}$ |
| $H_1: \mu > \mu_0$ | $\sqrt{-}(\bar{\mathbf{v}}_{-}, \bar{\mathbf{v}}_{-})/C$ | TC < | D(T < I) |
| | $\sqrt{n}(\bar{X}-\mu_0)/S$ | TS < t | $P\{T_{n-1} \le t\}$ |
| $H_1: \mu < \mu_0$ | | $-t_{\alpha,n-1}$ | |

 T_{n-1} is a t-random variable with (n - 1) degrees of freedom: $P(T_{n-1} > t_{\alpha,n-1}) = \alpha$.

Problem

A public health official claims that the mean home water use is at most 350 gallons a day. To verify this claim, a study of 20 randomly selected homes was instigated with the result that the average daily water uses of these 20 homes were as follows: 340 344 362 375 356 386 354 364 332 402 340 355 362 322 372 324 318 360 338 370 Do the data contradict the official's claim?

```
import numpy as np
from scipy import stats
# Given data
data = [340, 344, 362, 375, 356, 386, 354, 364, 332, 402, 340, 355, 362, 322, 372, 324, 318, 360, 338, 370]
sample mean = np.mean(data)
sample\_std = np.std(data, ddof=1)
sample\_size = len(data)
population mean = 350
alpha = 0.05
# Calculate the T-score
t_score = (sample_mean - population_mean) / (sample_std / np.sqrt(sample_size))
# Calculate the p-value
p_value = 2 * (1 - stats.t.cdf(np.abs(t_score), df=sample_size-1))
# Determine if we reject the null hypothesis
reject_null = p_value < alpha
# Output the results
print(f"Sample Mean: {sample mean}")
print(f"Sample Standard Deviation: {sample_std}")
print(f"T-score: {t_score}")
print(f"P-value: {p_value}")
print(f"Reject the null hypothesis: {reject_null}")
```

Sample Mean: 353.8

Sample Standard Deviation: 21.847798877449275

T-score: 0.7778411328447066 P-value: 0.4462410900531899 Reject the null hypothesis: False

5 Hypothesis Tests Concerning the variance of a normal population

- Let X_1, X_2, \dots, X_n be a sample of size n from a normal distribution with an unknown mean μ and variance σ^2 .
- We are interested in testing the null hypothesis:

$$H_0: \sigma^2 = \sigma_0^2$$

• Against the alternative hypothesis:

$$H_1: \sigma^2 \neq \sigma_0^2$$

- Where σ_0^2 is a specified constant.
- We know from the discussion on sampling distribution, $\frac{(n-1)S^2}{\sigma_0^2}$ has a chi-squared distribution with (n-1) degrees of freedom.

$$\frac{(n-1)S^2}{\sigma_0^2} \sim \chi_{n-1}^2$$

Also,

$$P_{H_0}\left(\chi^2_{1-\alpha/2,\;n-1} \leq \frac{(n-1)S^2}{\sigma_0^2} \leq \chi^2_{\alpha/2,\;n-1}\right) = 1 - \alpha$$

• In this case the test statistic (TS) is:

$$TS = \frac{(n-1)S^2}{\sigma_0^2}$$

• The p-value for this case is:

$$p-value = 2 \ \min \left(P(\chi^2_{n-1} < TS), 1 - P(\chi^2_{n-1} < TS) \right)$$

- 5 Hypothesis Tests Concerning the variance of a normal population
- 5.0.1 Hypothesis Testing Summary: chi-square Test

Table 5.1: Caption: Summary of hypothesis testing for a sample from a $N(\mu, \sigma^2)$ population.

| рораналон. | |
|-----------------------|--|
| Sample and Population | Details |
| Sample | $\{X_1, X_2,, X_n\}$ |
| Population | $N(\mu,\sigma^2)$ |
| Sample Mean | $\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$ |
| Sample Variance | $S^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (X_i - \bar{X})^2$ |
| Significance Level | α |

| Hypothesis | Test Statistic (TS) | Reject if | p-Value if $TS = t$ |
|---|-------------------------------|---|--|
| $H_0: \sigma^2 = \sigma_0^2$ vs $H_1: \sigma^2 \neq \sigma_0^2$ | $\frac{(n-1)S^2}{\sigma_0^2}$ | $TS \notin \left[\chi^2_{1-\alpha/2, n-1}, \right]$ | $2 \min_{\substack{\chi^2_{\alpha/2, n-1}}} (P(\chi^2_{n-1} < t), 1 - P(\chi^2_{n-1} < t)) $ |
| $H_0: \sigma^2 \le \sigma_0^2$ vs | $\frac{(n-1)S^2}{\sigma_0^2}$ | $TS > \chi^2_{\alpha, n-1}$ | $P\left(\chi_{n-1}^2 \ge t\right)$ |
| $\begin{aligned} H_1: \sigma^2 &> \sigma_0^2 \\ H_0: \sigma^2 &\geq \sigma_0^2 \\ \text{vs} \\ H_1: \sigma^2 &< \sigma_0^2 \end{aligned}$ | $\frac{(n-1)S^2}{\sigma_0^2}$ | $TS < \chi^2_{1-\alpha, n-1}$ | $P\left(\chi_{n-1}^2 \le t\right)$ |

Problem

A machine that automatically controls the amount of ribbon on a tape has recently been installed. This machine will be judged to be effective if the standard deviation sigma of the amount of ribbon on a tape is less than .15 cm. If a sample of 20 tapes yields a sample variance of $S^2 = .025 \text{ cm}^2$, are we justified in concluding that the machine is ineffective? Assume the level of significance as 0.05.

```
import numpy as np
from scipy.stats import chi2

# Given data
sample_variance = 0.025
sample_size = 20
population_variance = 0.15**2
alpha = 0.05
```

```
# Calculate the test statistic
test_statistic = (sample_size - 1) * sample_variance / population_variance

# Calculate the critical values
chi2_critical_low = chi2.ppf(alpha / 2, df=sample_size - 1)
chi2_critical_high = chi2.ppf(1 - alpha / 2, df=sample_size - 1)

# Determine if we reject the null hypothesis
reject_null = test_statistic < chi2_critical_low or test_statistic > chi2_critical_high

# Output the results
print(f"Test Statistic: {test_statistic}")
print(f"Chi-square Critical Low: {chi2_critical_low}")
print(f"Chi-square Critical High: {chi2_critical_high}")
print(f"Reject the null hypothesis: {reject_null}")
```

Test Statistic: 21.1111111111111114

Chi-square Critical Low: 8.906516481987971 Chi-square Critical High: 32.85232686172969

Reject the null hypothesis: False

References

Ross, Sheldon. 2009. "Probability and Statistics for Engineers and Scientists." Elsevier, New Delhi 16: 32-33.