

# Northeastern University

## College of Professional Studies

### Chi-Square testing and ANOVA

#### Overview and Rationale

In this assignment, you will use your knowledge of chi-square and ANOVA testing to solve various types of problems.

#### Course Outcomes

This assignment is directly linked to the following key learning outcomes from the course syllabus:

- CLO4: Test a distribution for goodness of fit, using chi-square.
- CLO5: Test two variables for independence, using chi-square.
- CLO6: Test proportions for homogeneity, using chi-square.
- CLO7: Use the one-way ANOVA technique to determine if there is a significant difference among three or more means.
- CLO9: Use the two-way ANOVA technique to determine if there is a significant difference in the main effects or interaction.

#### Assignment Summary

Complete the following problems using R and/or MS Excel. Be sure to show your work and include the hypothesis tests, the critical values, the computed test values, and the resulting decisions where applicable.

#### Section 11-1

*Perform these steps:*

- a. *State the hypotheses and identify the claim.*
- b. *Find the critical value.*
- c. *Compute the test value.*
- d. *Make the decision.*
- e. *Summarize the results.*

*Use the traditional method of hypothesis testing unless otherwise specified. Assume all assumptions are met.*

**6. Blood Types** A medical researcher wishes to see if hospital patients in a large hospital have the same blood type distribution as those in the general population. The distribution for the general population is as follows: type A, 20%; type B, 28%; type O, 36%; and type AB = 16%. He selects a random sample of 50 patients and finds the following: 12 have type A blood, 8 have type B, 24 have type O, and 6 have type AB blood.

At  $\alpha = 0.10$ , can it be concluded that the distribution is the same as that of the general population?

**8. On-Time Performance by Airlines** According to the Bureau of Transportation Statistics, on-time performance by the airlines is described as follows:

| Action                         | % of Time |
|--------------------------------|-----------|
| On time                        | 70.8      |
| National Aviation System delay | 8.2       |
| Aircraft arriving late         | 9.0       |

# Northeastern University

## College of Professional Studies

Other (because of weather and other conditions)

12.0

Records of 200 randomly selected flights for a major airline company showed that 125 planes were on time; 40 were delayed because of weather, 10 because of a National Aviation System delay, and the rest because of arriving late. At  $\alpha = 0.05$ , do these results differ from the government's statistics?

Source: [www.transtats.bts.gov](http://www.transtats.bts.gov)

### Section 11-2

Perform the following steps.

- State the hypotheses and identify the claim.
- Find the critical value.
- Compute the test value.
- Make the decision.
- Summarize the results.

Use the traditional method of hypothesis testing unless otherwise specified. Assume all assumptions are valid.

**8. Ethnicity and Movie Admissions** Are movie admissions related to ethnicity? A 2014 study indicated the following numbers of admissions (in thousands) for two different years. At the 0.05 level of significance, can it be concluded that movie attendance by year was dependent upon ethnicity?

|      | Caucasian | Hispanic | African American | Other |
|------|-----------|----------|------------------|-------|
| 2013 | 724       | 335      | 174              | 107   |
| 2014 | 370       | 292      | 152              | 140   |

Source: MPAA Study.

**10. Women in the Military** This table lists the numbers of officers and enlisted personnel for women in the military. At  $\alpha = 0.05$ , is there sufficient evidence to conclude that a relationship exists between rank and branch of the Armed Forces?

|              | Officers | Enlisted |
|--------------|----------|----------|
| Army         | 10,791   | 62,491   |
| Navy         | 7,816    | 42,750   |
| Marine Corps | 932      | 9,525    |
| Air Force    | 11,819   | 54,344   |

Source: New York Times Almanac.

### Section 12-1

Assume that all variables are normally distributed, that the samples are independent, that the population variances are equal, and that the samples are simple random samples, one from each of the populations. Also, for each exercise, perform the following steps.

- State the hypotheses and identify the claim.
- Find the critical value.
- Compute the test value.
- Make the decision.

# Northeastern University

## College of Professional Studies

*e. Summarize the results,*

*and explain where the differences in the means are.*

*Use the traditional method of hypothesis testing unless otherwise specified.*

**8. Sodium Contents of Foods** The amount of sodium (in milligrams) in one serving for a random sample of three different kinds of foods is listed. At the 0.05 level of significance, is there sufficient evidence to conclude that a difference in mean sodium amounts exists among condiments, cereals, and desserts?

| Condiments | Cereals | Desserts |
|------------|---------|----------|
| 270        | 260     | 100      |
| 130        | 220     | 180      |
| 230        | 290     | 250      |
| 180        | 290     | 250      |
| 80         | 200     | 300      |
| 70         | 320     | 360      |
| 200        | 140     | 300      |
|            |         | 160      |

*Source: The Doctor's Pocket Calorie, Fat, and Carbohydrate Counter.*

### Section 12-2

*Perform a complete one-way ANOVA. If the null hypothesis is rejected, use either the Scheffé or Tukey test to see if there is a significant difference in the pairs of means. Assume all assumptions are met.*

**10. Sales for Leading Companies** The sales in millions of dollars for a year of a sample of leading companies are shown. At  $\alpha = 0.01$ , is there a significant difference in the means?

|        | Chocolate |        |
|--------|-----------|--------|
| Cereal | Candy     | Coffee |
| 578    | 311       | 261    |
| 320    | 106       | 185    |
| 264    | 109       | 302    |
| 249    | 125       | 689    |
| 237    | 173       |        |

*Source: Information Resources, Inc.*

**12. Per-Pupil Expenditures** The expenditures (in dollars) per pupil for states in three sections of the country are listed. Using  $\alpha = 0.05$ , can you conclude that there is a difference in means?

| Eastern third | Middle third | Western third |
|---------------|--------------|---------------|
| 4946          | 6149         | 5282          |
| 5953          | 7451         | 8605          |
| 6202          | 6000         | 6528          |
| 7243          | 6479         | 6911          |

# Northeastern University

## College of Professional Studies

6113

Source: New York Times Almanac.

### Section 12-3

Assume that all variables are normally or approximately normally distributed, that the samples are independent, and that the population variances are equal.

- State the hypotheses.
- Find the critical value for each F test.
- Complete the summary table and find the test value.
- Make the decision.
- Summarize the results. (Draw a graph of the cell means if necessary.)

**10. Increasing Plant Growth** A gardening company is testing new ways to improve plant growth. Twelve plants are randomly selected and exposed to a combination of two factors, a “Grow-light” in two different strengths and a plant food supplement with different mineral supplements. After a number of days, the plants are measured for growth, and the results (in inches) are put into the appropriate boxes.

|              | Grow-light    |               |
|--------------|---------------|---------------|
|              | 1             | 2             |
| Plant food A | 9.2, 9.4, 8.9 | 8.5, 9.2, 8.9 |
| Plant food B | 7.1, 7.2, 8.5 | 5.5, 5.8, 7.6 |

Can an interaction between the two factors be concluded? Is there a difference in mean growth with respect to light? With respect to plant food? Use  $\alpha = 0.05$ .

Use R to complete the following steps. Be sure to include all code in an appendix at the end of your submission. Assume the expected frequencies are equal and  $\alpha = 0.05$ .

- Download the file ‘baseball.csv’ from the course resources and import the file into R.
- Perform EDA on the imported data set. Write a paragraph or two to describe the data set using descriptive statistics and plots. Are there any trends or anything of interest to discuss?
- Assuming the expected frequencies are equal, perform a Chi-Square Goodness-of-Fit test to determine if there is a difference in the number of wins by decade. Be sure to include the following:
  - State the hypotheses and identify the claim.
  - Find the critical value ( $\alpha = 0.05$ ) (From table in the book).
  - Compute the test value.
  - Make the decision. Clearly state if the null hypothesis should or should not be rejected and why.
  - Does comparing the critical value with the test value provide the same result as comparing the p-value from R with the significance level?

Here is some code to get you started. Be sure to import the dplyr and tidyverse packages.

```
# Extract decade from year
bb$Decade <- bb$Year - (bb$Year %% 10)

# Create a wins table by summing the wins by decade
wins <- bb %>%
  group_by(Decade) %>%
```

# Northeastern University

## College of Professional Studies

summarize(wins = sum(W))

```
%>%  
as.tibble()
```

4. Download the file 'crop\_data.csv' from the course resources and import the file into R.
5. Perform a Two-way ANOVA test using *yield* as the dependent variable and *fertilizer* and *density* as the independent variables. Explain the results of the test. Is there reason to believe that fertilizer and density have an impact on yield?

\*\* Be sure to convert the variables density, fertilizer and block to R factors.  
\*\* Include a null and alternate hypothesis for both factors and the interaction.

### Report

Refer to the attached rubric for more details on the report. The report should contain a well written cover/title page, introduction, body, conclusion, and references. It must follow APA format and have at least 1000 words (excluding title page and references page. All R code used for your report should be included in an appendix at the end of the report.

Graphs, figures, charts, and tables are very useful visual effects to communicate your results and impress your readers. However, such items should not be included in the report unless they are well described and interpreted. Please use subtitles to make your assignment more reader friendly as well.

### Format & Guidelines

The report should follow the following format:

- (i) Title page
- (ii) Introduction
- (iii) Analysis
- (iv) Conclusion/Interpretations
- (v) References

### Assignment Rubric

| Category                          | Above Standards  | Meets Standards   | Approaching Standards  | Below Standards   |
|-----------------------------------|--|---|--|---|
| <b>Introduction</b><br><b>15%</b> | Clearly and briefly introduces the goals of the project, the question that needs to be answered and the methods used in the analysis. The goals, questions and methods outlined are consistent with one another. | Introduction provides a brief and intelligible overview of the goals and methods of the assignment. | Introduction provides an overview of the goals and methods of the assignment, but is ambiguous or not concise. | Does not introduce project goals, project questions or methods. |

# Northeastern University

## College of Professional Studies

|   |   |   |   |  |
|---|---|---|---|--|
| <b>Analysis</b><br><br>25%  | Incorporates R code and the outputs. Provides detailed analysis of the output focusing on significance results. Uses visualizations to make major points.   | Provides all R code and the outputs. Includes interpretation of the output, graphs, figures, charts, and tables and the significance of the results in the analysis.                        | Provides R codes and outputs, but the R code does not match the outputs or is missing some code or outputs. Includes limited interpretations, charts, and tables and the significance of the results in the analysis. | Does not provide R code or its outputs or minimal R code is provided. Includes few interpretations, charts, or tables. Does not identify the significance of the results in the analysis.      |
| <b>Data Visualizations</b><br><br>25%                                     | Data visualizations are appropriate for the level and type of analysis. . Uses graphs, figures, charts, and tables to increase visual effects of the main points being made based on the results. | Data visualizations are appropriate for the level and type of analysis. Graphs, figures and tables communicate insights and significance to the reader.                                     | Data visualization are useful for the level and type of analysis, but graphs, figures and tables do not clearly communicate the significance of the results to the reader.  | Data visualization are used minimally or not at all. If graphs, figures and tables are used, it is unclear what they are intended to communicate or why.                                       |
| <b>Interpretation &amp; Conclusions</b><br><br>25%                        | Wraps up the findings in a conclusion that provides an answer to the question(s) posed in the introduction. Makes specific recommendations based on the data presented.                           | The conclusion summarizes and makes sense of the results, making good points that reflect clear understanding of the assignment material.   | The conclusion summarizes and makes sense of the results, making good points that reflect a basic understanding of the assignment material.   | The conclusion does not summarize or attempt to make sense of the results. Conclusions do not reflect an understanding or reflect a misunderstanding of the material.                          |
| <b>Report: Writing Mechanics, Title Page, &amp; References</b><br><br>10% | There are no noticeable errors in grammar, spelling, and punctuation; and completely correct usage of title page, citations, and references. The report contains approximately 1,000 words.       | There are no noticeable errors in grammar, spelling, and punctuation; and completely correct usage of title page, citations, and references. The report contains approximately 1,000 words. | There are very few errors in grammar, spelling, and punctuation; and completely correct usage of title page, citations, and references. The report contains approximately 1,000 words.                                | There are more than five errors in grammar, spelling, and punctuation; or the usage of title page, citations, and references are incomplete; or the report contains far less than 1,000 words. |

