

Cookbook Data Analysis with Stata and R

Manuel Oliveira

2024-08-01

Table of contents

Preface	5
1 Introduction	6
2 Chapter 1: Getting Started	7
2.1 Introduction	7
2.2 Installing R	7
2.2.1 Windows	7
2.2.2 Mac	7
2.3 Installing Stata	8
2.3.1 Windows	8
2.3.2 Mac	8
2.4 Setting Up Your Environment	9
2.4.1 R Setup	9
2.4.2 Stata Setup	9
2.5 Verification	9
2.5.1 R	9
2.5.2 Stata	10
3 t-test	11
3.1 Brief Explanation	11
3.2 Research Question and Hypothesis	11
3.3 R packages	11
3.4 Simulated Dataset	13
3.4.1 In R	13
3.4.2 In Stata	13
3.5 Statistical Analysis	14
3.5.1 Inspecting Data Descriptives and Plotting	14
3.5.2 T-test	15
3.6 Explanation of Relevant Terms	16
3.7 Interpretation Questions	16
3.8 Comparison Tables	17
3.8.1 Assumptions and Statistical Analysis Commands	17
3.8.2 Simulated Dataset Generation	17

4	ANOVA	18
4.1	Introduction	18
4.2	Example Question	18
4.3	Required Packages (R)	18
4.4	Simulating the Dataset in R	19
4.5	Simulating the Dataset in Stata	20
4.6	Visualizing the Descriptives in R	20
4.7	Visualizing the Descriptives in Stata	21
4.8	Running the ANOVA in R	21
4.9	Running the ANOVA in Stata	22
4.10	Interpreting the Output	22
	4.10.1 In R	22
	4.10.2 In Stata	22
4.11	Post-hoc Testing in R	22
4.12	Post-hoc Testing in Stata	23
4.13	Plotting the Results in R	23
4.14	Plotting the Results in Stata	25
4.15	Assumptions	25
4.16	Syntax Comparison: R vs Stata	25
5	Linear Regression	27
5.1	Introduction	27
5.2	Example Question	27
5.3	Dataset Simulation in R	27
5.4	Dataset Simulation in Stata	28
5.5	Performing Linear Regression	28
	5.5.1 R	28
	5.5.2 Stata	28
5.6	Assumptions	28
6	Multilevel Regression	29
6.1	Introduction	29
6.2	Example Question	29
6.3	Dataset Simulation in R	29
6.4	Dataset Simulation in Stata	30
6.5	Performing Multilevel Regression	30
	6.5.1 R	30
	6.5.2 Stata	30
6.6	Assumptions	31
7	Logistic Regression	32
7.1	Introduction	32
7.2	Example Question	32

7.3	Required Packages (R)	32
7.4	Simulating the Dataset in R	33
7.5	Simulating the Dataset in Stata	33
7.6	Visualizing the Descriptives in R	34
7.7	Visualizing the Descriptives in Stata	35
7.8	Running the Logistic Regression in R	36
7.8.1	Output interpretation	37
7.9	Running the Logistic Regression in Stata	37
7.9.1	Output description	38
7.10	Plotting the Results in R	38
7.11	Plotting the Results in Stata	39
7.12	Assumptions	40
7.13	R	40
7.13.1	Binary outcome	40
7.13.2	Independence	41
7.13.3	No Multicollinearity	42
7.13.4	Large sample size	42
7.14	Stata	43
7.14.1	Binary outcome	43
7.14.2	Independence	43
7.14.3	Linearity of logit	43
7.14.4	No Multicollinearity	43
7.14.5	Large sample size	44
7.15	Syntax Comparison: R vs Stata	44
8	Summary	45
	References	46

Preface

This is a Quarto book.

To learn more about Quarto books visit <https://quarto.org/docs/books>.

```
cat(" This is a test again")
```

This is a test again

1 Introduction

This is a book created from markdown and executable code.

See Knuth (1984) for additional discussion of literate programming.

```
1 + 1
```

```
[1] 2
```

2 Chapter 1: Getting Started

2.1 Introduction

This chapter provides a quick tutorial on how to install and set up R and Stata on both Windows and Mac computers. By the end of this chapter, you'll have the necessary tools ready to begin your analysis.

2.2 Installing R

2.2.1 Windows

1. **Download R:**

- Go to the [R Project website](#).
- Click on “Download R for Windows.”
- Click on “base” to download the base R package.

2. **Install R:**

- Run the downloaded `.exe` file.
- Follow the installation instructions, accepting the default settings.

3. **Install RStudio** (Optional but recommended):

- Download RStudio from the [RStudio website](#).
- Run the installer and follow the setup instructions.

2.2.2 Mac

1. **Download R:**

- Visit the [R Project website](#).
- Click on “Download R for macOS.”

2. **Install R:**

- Open the downloaded `.pkg` file.

- Follow the installation instructions.
3. **Install RStudio** (Optional but recommended):
 - Download RStudio from the [RStudio website](#).
 - Open the `.dmg` file and drag RStudio to your Applications folder.

2.3 Installing Stata

2.3.1 Windows

1. **Obtain a License:**
 - Stata is commercial software. Ensure you have a valid license.
2. **Download Stata:**
 - Go to the [Stata website](#) and log in to your account to download the installer.
3. **Install Stata:**
 - Run the downloaded `.exe` file.
 - Follow the installation instructions, entering your license information when prompted.

2.3.2 Mac

1. **Obtain a License:**
 - Make sure you have a valid license for Stata.
2. **Download Stata:**
 - Visit the [Stata website](#) and log in to your account to download the installer.
3. **Install Stata:**
 - Open the downloaded `.dmg` file.
 - Drag the Stata application to your Applications folder.
 - Launch Stata and enter your license information.

2.4 Setting Up Your Environment

2.4.1 R Setup

1. **Open RStudio** (or R GUI if not using RStudio).

2. **Install Essential Packages:**

- Open the Console and run:

```
install.packages(c("tidyverse", "lme4", "ggplot2"))
```

3. **Create a New Project** (Optional but recommended in RStudio):

- Go to “File” > “New Project” > “New Directory” > “New Project.”
- Choose a location and name for your project, then click “Create Project.”

2.4.2 Stata Setup

1. **Open Stata.**

2. **Set a Working Directory:**

- Use the command:

```
cd "path/to/your/directory"
```

Replace "path/to/your/directory" with the path where you want to save your files.

3. **Creating Do-Files:**

- Go to “File” > “New Do-file Editor.”
- Save the Do-file in your working directory.

2.5 Verification

2.5.1 R

1. **Test Installation:**

- In RStudio or R GUI, type:

```
print("R is working!")
```

- If you see the output `[1] "R is working!"`, your installation is successful.

2. Load a Package:

- Run:

```
library(ggplot2)
print("ggplot2 is loaded!")
```

2.5.2 Stata

1. Test Installation:

- In the Command window, type:

```
display "Stata is working!"
```

- If you see the output `Stata is working!`, your installation is successful.

2. Check Version:

- Type:

```
about
```

- This will display the version of Stata installed.

With your environment set up, you're now ready to start performing analyses using R and Stata!

3 t-test

3.1 Brief Explanation

The t-test, proposed by William Sealy Gosset under the pseudonym “Student” in 1908, is used to determine if there is a significant difference between the means of two groups. It is commonly used when the sample sizes are small and the population variance is unknown.

Statistic	Description
Proposed by	William Sealy Gosset (1908)
Purpose	Compare means of two groups
When to use	Small sample sizes, unknown population variance
Example question	Is there a significant difference in user satisfaction between two versions of a software?
Analytical goal	Determine if the mean satisfaction scores differ significantly between the two versions

3.2 Research Question and Hypothesis

Research Question: Does the new interface design improve user satisfaction compared to the old design?

Hypothesis: Users will report higher satisfaction scores with the new interface design compared to the old design.

3.3 R packages

```
# Load the necessary packages
library(tidyverse) # used for data manipulation and visualization
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2     3.5.0      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.1
v purrr       1.0.2
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(gggrain) # used for raincloud plots
```

Registered S3 methods overwritten by 'ggpp':

```
method          from
heightDetails.titleGrob ggplot2
widthDetails.titleGrob  ggplot2
```

```
library(cowplot) # for cowplot theme in ggplot
```

Attaching package: 'cowplot'

The following object is masked from 'package:lubridate':

```
stamp
```

```
library(Statamarkdown) # to run Stata commands in an R environment
```

Stata found at C:/Program Files/Stata18/StataBE-64.exe

The 'stata' engine is ready to use.

```
# Statamarkdown configuration
stataexe <- "C:/Program Files/Stata18/StataBE-64.exe" # Add your own path to the Stata executable
knitr::opts_chunk$set(engine.path=list(stata=stataexe))

# to install any missing packages go to the Terminal and run the command: install.packages("I
```

3.4 Simulated Dataset

3.4.1 In R

```
set.seed(123)
n <- 30
old_design <- rnorm(n, mean = 70, sd = 10)
new_design <- rnorm(n, mean = 75, sd = 10)
data <- data.frame(
  group = rep(c("Old Design", "New Design"), each = n),
  satisfaction = c(old_design, new_design)
)
write.csv(data, "satisfaction_data.csv", row.names = FALSE)
```

3.4.2 In Stata

```
clear
set seed 123
set obs 30
gen group = "Old Design"
gen satisfaction = rnormal(70, 10)
save old_design.dta, replace
```

Number of observations (_N) was 0, now 30.

file old_design.dta saved

```
clear
set obs 30
gen group = "New Design"
gen satisfaction = rnormal(75, 10)
save new_design.dta, replace
```

Number of observations (_N) was 0, now 30.

```
file new_design.dta saved
```

```
use old_design.dta
append using new_design.dta
save satisfaction_data.dta, replace
```

```
file satisfaction_data.dta saved
```

3.5 Statistical Analysis

3.5.1 Inspecting Data Descriptives and Plotting

3.5.1.1 In Stata

```
use satisfaction_data.dta
summarize satisfaction
graph box satisfaction, over(group)
```

Variable	Obs	Mean	Std. dev.	Min	Max
satisfaction	60	72.69432	11.03765	50.47923	99.36316

3.5.1.2 In R

```
data <- read.csv("satisfaction_data.csv")
summary(data$satisfaction)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
50.33	65.48	73.26	73.16	80.62	96.69

```
p <- ggplot(data, aes(x = group, y = satisfaction)) +
  geom_boxplot() +
  theme_cowplot() +
  labs(title = "Satisfaction Scores by Group", x = "Group", y = "Satisfaction Score")
p
```

A box plot comparing the Satisfaction Score for two groups: New Design and Old Design. The Y-axis represents the Satisfaction Score, ranging from 50 to 90. The X-axis is labeled 'Group'.

The New Design group shows a median satisfaction score of approximately 75.5, with a box spanning from about 72 to 82.5. The Old Design group shows a median satisfaction score of approximately 69, with a box spanning from about 63 to 75.

Both groups have whiskers extending to the minimum and maximum values. The New Design group has a minimum of approximately 59.5 and a maximum of approximately 96.5. The Old Design group has a minimum of approximately 50.5 and a maximum of approximately 88.

Group	Min	Q1	Median	Q3	Max
New Design	59.5	72.0	75.5	82.5	96.5
Old Design	50.5	63.0	69.0	75.0	88.0

3.5.2.1 In Stata

Two-sample t test with equal variances

```
diff = mean(New Desi) - mean(Old Desi)          t = 1.7059
H0: diff = 0                                     Degrees of freedom = 58
```

Ha: diff < 0
Pr(T < t) = 0.9533

Ha: diff != 0
Pr(|T| > |t|) = 0.0934

Ha: diff > 0
Pr(T > t) = 0.0467

3.5.2.2 In R

```
t.test(satisfaction ~ group, data = data)
```

Welch Two Sample t-test

```
data: satisfaction by group
t = 3.0841, df = 56.559, p-value = 0.003156
alternative hypothesis: true difference in means between group New Design and group Old Design
95 percent confidence interval:
 2.543416 11.965426
sample estimates:
mean in group New Design mean in group Old Design
      76.78338             69.52896
```

3.6 Explanation of Relevant Terms

Term	Description
t-value	The calculated difference represented in units of standard error
p-value	The probability that the observed difference occurred by chance
Confidence Interval	The range within which the true population mean difference lies with a certain level of confidence

3.7 Interpretation Questions

1. What does a significant p-value indicate in the context of this t-test?
2. How would you interpret the confidence interval in this analysis?

Solutions:

1. A significant p-value indicates that there is a statistically significant difference between the satisfaction scores of the two groups.
2. The confidence interval provides a range of values within which we can be confident that the true mean difference lies.

3.8 Comparison Tables

3.8.1 Assumptions and Statistical Analysis Commands

Step	R Command	Stata Command
Descriptive Statistics	<code>summary(data\$satisfaction)</code>	<code>summarize satisfaction</code>
Box Plot	<code>ggplot(data, aes(x = group, y = satisfaction)) + geom_boxplot()</code>	<code>graph box satisfaction, over(group)</code>
T-Test	<code>t.test(satisfaction ~ group, data = data)</code>	<code>ttest satisfaction, by(group)</code>

3.8.2 Simulated Dataset Generation

Step	R Command	Stata Command
Generate Data	<code>rnorm(n, mean, sd)</code>	<code>rnormal(mean, sd)</code>
Save Data	<code>write.csv(data, "satisfaction_data.csv")</code>	<code>save satisfaction_data.dta, replace</code>

4 ANOVA

4.1 Introduction

This chapter covers ANOVA (Analysis of Variance), used to compare the means across multiple groups. We will use an example dataset to investigate whether the design of a user interface (UI) affects the time users spend on a website.

4.2 Example Question

Does the design of a user interface (UI) influence the time users spend on a website?

4.3 Required Packages (R)

```
# Load the necessary packages
library(tidyverse) # used for data manipulation and visualization

-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.0      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.1
v purrr      1.0.2

-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(car) # provides tools for ANOVA and regression diagnostics
```

Loading required package: carData

Attaching package: 'car'

The following object is masked from 'package:dplyr':

recode

The following object is masked from 'package:purrr':

some

```
# to install any missing packages go to the Terminal and run the command: install.packages("l
```

4.4 Simulating the Dataset in R

```
# Setting a seed for reproducibility
set.seed(123)

# Simulating data
n_groups <- 3 # Number of UI designs
n_per_group <- 50 # Number of users per group

# Creating a factor variable for UI design
ui_design <- factor(rep(1:n_groups, each = n_per_group))

# Simulating time spent data with different means for each UI design
time_spent <- rnorm(n_groups * n_per_group, mean = rep(c(20, 25, 22), each = n_per_group), s

# Creating a data frame
data <- data.frame(ui_design, time_spent)

# Viewing the first few rows of the dataset
head(data)
```

	ui_design	time_spent
1	1	17.19762
2	1	18.84911
3	1	27.79354

4	1	20.35254
5	1	20.64644
6	1	28.57532

4.5 Simulating the Dataset in Stata

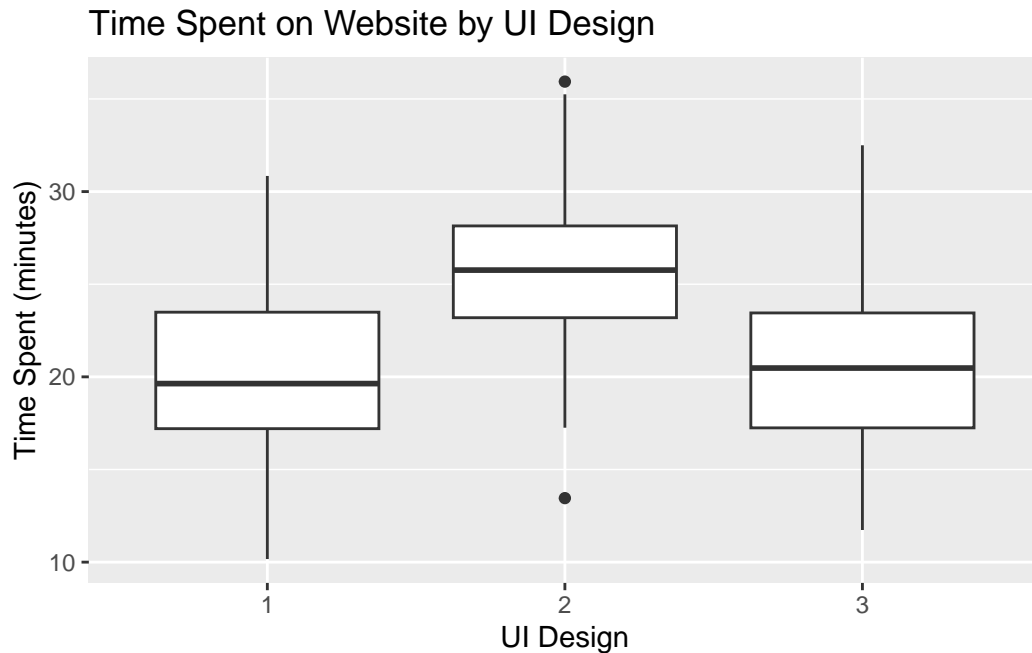
```
* Set seed for reproducibility
set seed 123

* Simulate data
set obs 150
gen ui_design = ceil(_n/50)
gen time_spent = rnormal(20 + (ui_design==2)*5 + (ui_design==3)*2, 5)

* View the first few rows
list in 1/10
```

4.6 Visualizing the Descriptives in R

```
# Plotting the distribution of time spent across different UI designs
ggplot(data, aes(x = ui_design, y = time_spent)) +
  geom_boxplot() +
  labs(title = "Time Spent on Website by UI Design",
       x = "UI Design",
       y = "Time Spent (minutes)")
```



4.7 Visualizing the Descriptives in Stata

```
* Box plot of time spent by UI design
graph box time_spent, over(ui_design) title("Time Spent on Website by UI Design") ///
    ytitle("Time Spent (minutes)") xtitle("UI Design")
```

4.8 Running the ANOVA in R

```
# Performing ANOVA
anova_model <- aov(time_spent ~ ui_design, data = data)

# Viewing the summary of the ANOVA model
summary(anova_model)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
ui_design	2	937	468.7	21.18	8.29e-09 ***
Residuals	147	3253	22.1		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

4.9 Running the ANOVA in Stata

```
* Perform ANOVA
anova time_spent ui_design
```

4.10 Interpreting the Output

4.10.1 In R

The ANOVA table provides the following key pieces of information: - **Df**: Degrees of freedom associated with the sources of variance. - **Sum Sq**: Sum of squares, which measures the total variation for each source. - **Mean Sq**: Mean square, calculated as Sum Sq divided by Df. - **F value**: The F-statistic, calculated as the ratio of mean square values. - **Pr(>F)**: The p-value associated with the F-statistic.

4.10.2 In Stata

The output of the ANOVA in Stata provides similar information: - **Source**: Lists the sources of variance. - **Partial SS**: Partial sum of squares for each source. - **df**: Degrees of freedom associated with each source. - **MS**: Mean square for each source, calculated as SS/df. - **F**: The F-statistic for each source. - **Prob > F**: The p-value associated with the F-statistic.

If the p-value is less than the significance level (typically 0.05), we reject the null hypothesis that all group means are equal.

4.11 Post-hoc Testing in R

```
# Performing Tukey's Honest Significant Difference test
tukey_test <- TukeyHSD(anova_model)

# Viewing the Tukey test results
tukey_test
```

Tukey multiple comparisons of means
95% family-wise confidence level

```
Fit: aov(formula = time_spent ~ ui_design, data = data)
```

```
$ui_design
      diff      lwr      upr    p adj
2-1  5.5600236  3.332272  7.787775 0.0000001
3-1  0.5584801 -1.669271  2.786231 0.8237890
3-2 -5.0015435 -7.229295 -2.773792 0.0000012
```

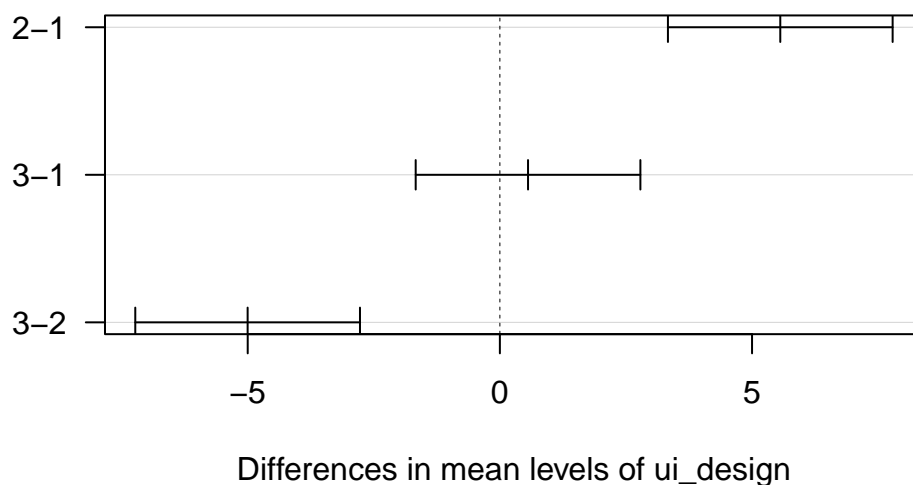
4.12 Post-hoc Testing in Stata

```
* Perform Bonferroni post-hoc test
oneway time_spent ui_design, bonferroni
```

4.13 Plotting the Results in R

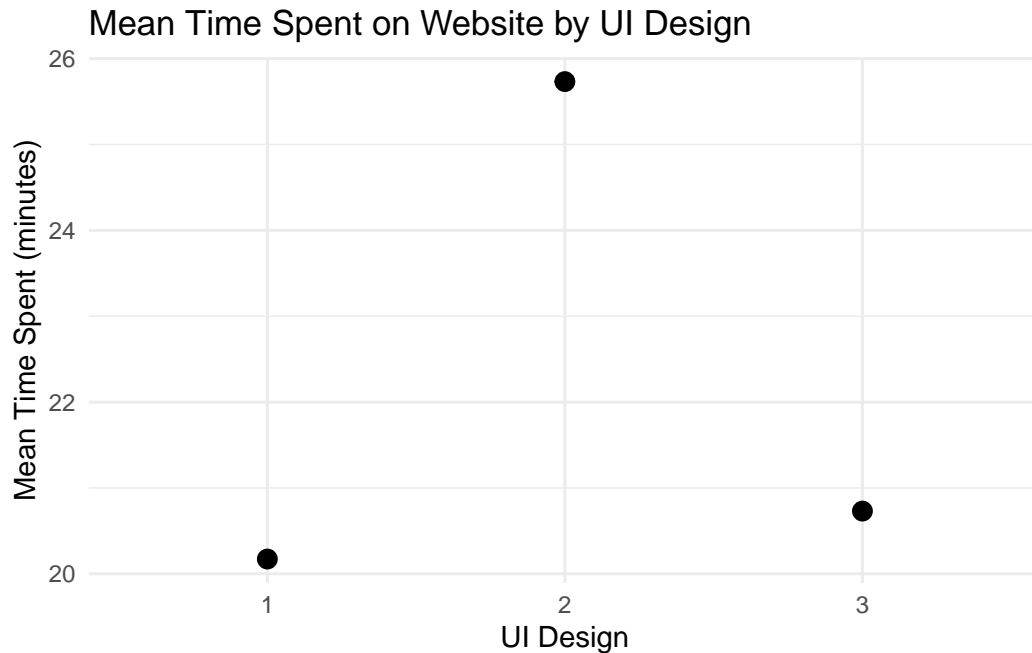
```
# Plotting the results of the Tukey HSD test
plot(tukey_test, las = 1)
```

95% family-wise confidence level



```
# Creating a plot to visualize group means with confidence intervals
ggplot(data, aes(x = ui_design, y = time_spent)) +
  stat_summary(fun.data = mean_cl_normal, geom = "errorbar", width = 0.2) +
  stat_summary(fun = mean, geom = "point", size = 3) +
  labs(title = "Mean Time Spent on Website by UI Design",
       x = "UI Design",
       y = "Mean Time Spent (minutes)") +
  theme_minimal()
```

Warning: Computation failed in `stat_summary()`.
Caused by error in `fun.data()`:
! The package "Hmisc" is required.



4.14 Plotting the Results in Stata

```
* Plot group means with confidence intervals  
means time_spent, over(ui_design) ci
```

4.15 Assumptions

- **Independence:** Observations should be independent of each other.
- **Normality:** The residuals of the model should be normally distributed.
- **Homoscedasticity:** Variances across the groups should be equal.
- **Random Sampling:** The data should be randomly sampled from the population.

These assumptions should be checked to ensure the validity of the ANOVA results.

4.16 Syntax Comparison: R vs Stata

This table summarizes the main differences between R and Stata in terms of syntax for performing ANOVA analyses.

Task	R Command	Stata Command
Simulating Data	<code>rnorm()</code> for simulating normal distribution	<code>rnormal()</code> for simulating normal distribution
Setting Seed for Reproducibility	<code>set.seed(123)</code>	<code>set seed 123</code>
Creating a Factor Variable	<code>factor()</code>	<code>gen variable</code> and <code>egen group</code>
Visualizing Descriptives	<code>ggplot()</code> with <code>geom_boxplot()</code>	<code>graph box</code>
Running ANOVA	<code>aov()</code> and <code>summary()</code>	<code>anova</code>
Post-hoc Testing	<code>TukeyHSD()</code>	<code>oneway</code> with <code>bonferroni</code> option
Plotting Group Means with Confidence Intervals	<code>ggplot()</code> with <code>stat_summary()</code>	<code>means</code> with <code>ci</code> option

5 Linear Regression

5.1 Introduction

This chapter covers how to perform linear regression to study the relationship between variables. We'll use an example dataset that simulates the relationship between study time and performance on an online learning platform.

5.2 Example Question

How does the amount of time spent on an e-learning platform (in hours) affect the test scores of users?

5.3 Dataset Simulation in R

```
# Load necessary package
set.seed(123)

# Simulate data
n <- 100
study_time <- rnorm(n, mean = 10, sd = 2) # Average 10 hours
test_score <- 50 + 5 * study_time + rnorm(n, mean = 0, sd = 5) # Linear relationship with s

# Create a data frame
data <- data.frame(study_time, test_score)

# View the first few rows
head(data)
```

5.4 Dataset Simulation in Stata

```
* Set seed for reproducibility
set seed 123

* Simulate data
set obs 100
gen study_time = rnormal(10, 2)
gen test_score = 50 + 5 * study_time + rnormal(0, 5)

* View the first few rows
list in 1/10
```

5.5 Performing Linear Regression

5.5.1 R

```
# Fit the linear regression model
model <- lm(test_score ~ study_time, data = data)

# View the summary
summary(model)
```

5.5.2 Stata

```
* Fit the linear regression model
regress test_score study_time
```

5.6 Assumptions

- **Linearity:** The relationship between the independent and dependent variable should be linear.
- **Independence:** Observations should be independent of each other.
- **Homoscedasticity:** The residuals should have constant variance at every level of the independent variable.
- **Normality:** The residuals should be normally distributed.

6 Multilevel Regression

6.1 Introduction

This chapter covers multilevel regression, where data is nested. We will explore how user satisfaction with a mobile app is affected by time spent on the app, considering that users are nested within different age groups.

6.2 Example Question

Does time spent on a mobile app influence user satisfaction, and does this effect differ across age groups?

6.3 Dataset Simulation in R

```
# Load necessary package
set.seed(123)

# Simulate data
n_groups <- 5 # Number of age groups
n_per_group <- 50 # Number of users per group

age_group <- factor(rep(1:n_groups, each = n_per_group))
time_spent <- rnorm(n_groups * n_per_group, mean = 30, sd = 10)
satisfaction <- 3 + 0.2 * time_spent + as.numeric(age_group) + rnorm(n_groups * n_per_group,

# Create a data frame
data <- data.frame(age_group, time_spent, satisfaction)

# View the first few rows
head(data)
```

6.4 Dataset Simulation in Stata

```
* Set seed for reproducibility
set seed 123

* Simulate data
set obs 250
gen group = ceil(_n/50) // Age group
gen time_spent = rnormal(30, 10)
gen satisfaction = 3 + 0.2 * time_spent + group + rnormal(0, 2)

* Convert group to a factor
egen group_factor = group(group)

* View the first few rows
list in 1/10
```

6.5 Performing Multilevel Regression

6.5.1 R

```
# Load necessary package
library(lme4)

# Fit the multilevel model
model <- lmer(satisfaction ~ time_spent + (1 | age_group), data = data)

# View the summary
summary(model)
```

6.5.2 Stata

```
* Fit the multilevel model
mixed satisfaction time_spent || group:
```

6.6 Assumptions

- **Normality of residuals:** The residuals at each level of the model should be normally distributed.
- **Linearity:** The relationship between predictors and the outcome should be linear at each level of the model.
- **Independence:** Observations within each group should be independent.
- **Homoscedasticity:** The variance of residuals should be consistent across all levels of the hierarchy.

7 Logistic Regression

7.1 Introduction

This chapter covers logistic regression, which is used when the outcome variable is binary. We will use an example dataset to investigate whether the frequency of technical support contact predicts whether a user continues to use a software product.

7.2 Example Question

Does the frequency of contacting technical support predict whether a user will continue using a software product?

7.3 Required Packages (R)

```
# Load the necessary packages
library(tidyverse) # used for data manipulation and visualization

-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.0      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.1
v purrr      1.0.2

-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(broom) # for tidying the model output, making it easier to work with
library(Statamarkdown) # to run Stata commands in an R environment
```


Stata found at C:/Program Files/Stata18/StataBE-64.exe
The 'stata' engine is ready to use.

```
# Statamarkdown configuration
stataexe <- "C:/Program Files/Stata18/StataBE-64.exe" # Add your own path to the Stata executable
knitr::opts_chunk$set(engine.path=list(stata=stataexe))

# to install any missing packages go to the Terminal and run the command: install.packages("r"
```

7.4 Simulating the Dataset in R

```
# Setting a seed for reproducibility
set.seed(123)

# Simulating data
n <- 200
support_contact <- rpois(n, lambda = 2) # Number of contacts with support
continued_use <- rbinom(n, size = 1, prob = 1 / (1 + exp(-(-1 + 0.5 * support_contact))))

# Creating a data frame
data <- data.frame(support_contact, continued_use)

# Viewing the first few rows of the dataset
head(data)
```

	support_contact	continued_use
1	1	0
2	3	0
3	2	1
4	4	1
5	4	1
6	0	1

7.5 Simulating the Dataset in Stata

```

* Set seed for reproducibility
set seed 123

* Simulate data
set obs 200
gen support_contact = rpoisson(2)
gen continued_use = rbinomial(1, 1 / (1 + exp(-(-1 + 0.5 * support_contact))))

* Save to data file
save logreg_data.dta

* View the first few rows
list in 1/10

```

Number of observations (_N) was 0, now 200.

```

file logreg_data.dta already exists
r(602);

```

```

r(602);

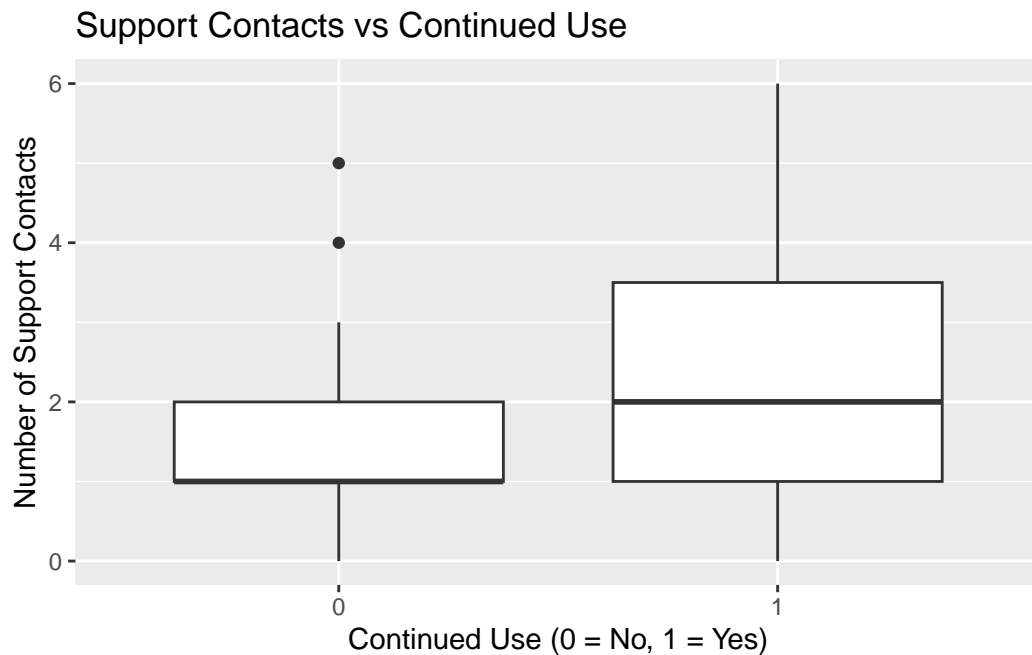
```

7.6 Visualizing the Descriptives in R

```

# Plotting the distribution of support contacts for users who continued vs those who didn't
ggplot(data, aes(x = factor(continued_use), y = support_contact)) +
  geom_boxplot() +
  labs(title = "Support Contacts vs Continued Use",
       x = "Continued Use (0 = No, 1 = Yes)",
       y = "Number of Support Contacts")

```



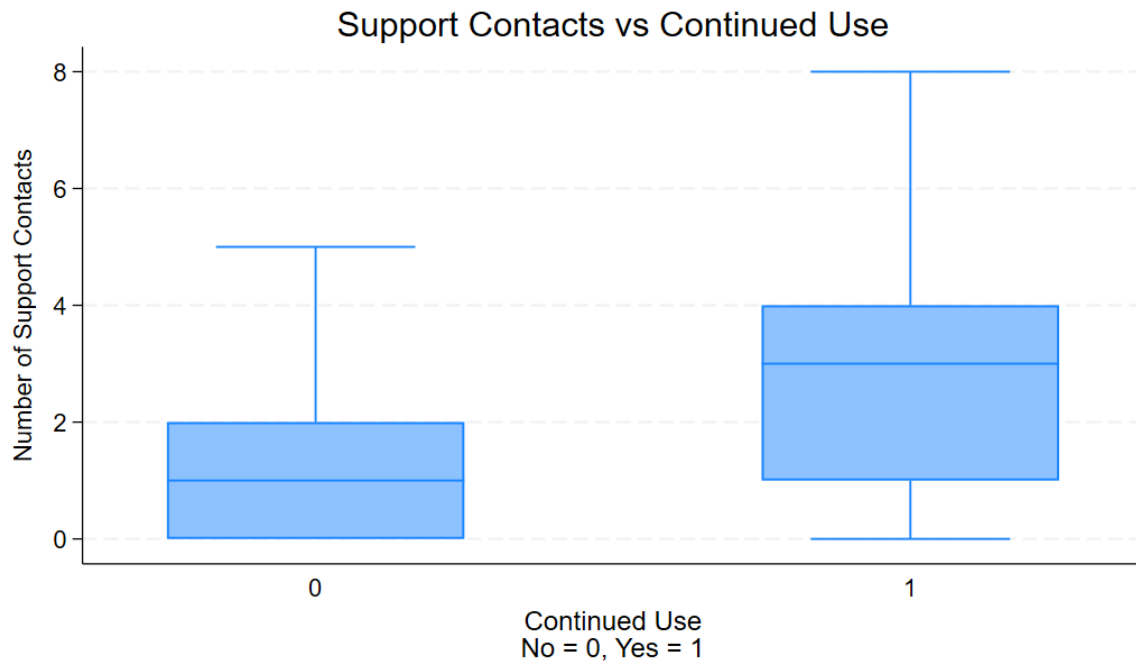
7.7 Visualizing the Descriptives in Stata

```
* NO NEED TO LOAD DATA AGAIN If USING STATA
use logreg_data.dta
```

```
* Data summary
summarize
```

Variable	Obs	Mean	Std. dev.	Min	Max
support_contact	200	1.875	1.523476	0	8
continued_use	200	.435	.4970011	0	1

```
* Box plot of support contacts by continued use
graph box support_contact, over(continued_use) title("Support Contacts vs Continued Use") b1
```



7.8 Running the Logistic Regression in R

```
# Fitting the logistic regression model
logistic_model <- glm(continued_use ~ support_contact, data = data, family = "binomial")

# Viewing the summary of the logistic regression model
summary(logistic_model)
```

Call:

```
glm(formula = continued_use ~ support_contact, family = "binomial",
    data = data)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-1.2023	0.3046	-3.947	7.90e-05	***
support_contact	0.7398	0.1453	5.092	3.54e-07	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 274.83 on 199 degrees of freedom
Residual deviance: 240.15 on 198 degrees of freedom
AIC: 244.15

Number of Fisher Scoring iterations: 4

7.8.1 Output interpretation

Term	Description
Coefficients	Estimates of the regression coefficients.
Std. Error	Standard errors of the coefficients.
z value	The test statistic for each coefficient.
Pr(> z)	The p-value associated with each coefficient, indicating whether it is statistically significant. If the p-value is less than the significance level (typically 0.05), we reject the null hypothesis that the coefficient is equal to zero.

7.9 Running the Logistic Regression in Stata

```
* Loading data to make it work in R environment, YOU DO NOT NEED TO LOEAD THE DATA AGAIN IN R  
use logreg_data.dta
```

```
* Fit the logistic regression model  
logit continued_use support_contact
```

> DATA AGAIN IN STATA IF YOU ALREADY LOADED IT BEFORE!

```
Iteration 0: Log likelihood = -136.93464  
Iteration 1: Log likelihood = -121.18683  
Iteration 2: Log likelihood = -121.17434  
Iteration 3: Log likelihood = -121.17434
```

Logistic regression

Number of obs = 200
LR chi2(1) = 31.52

Log likelihood = -121.17434

Prob > chi2 = 0.0000
Pseudo R2 = 0.1151

continued~e	Coefficient	Std. err.	z	P> z	[95% conf. interval]	
support_co~t	.589038	.1166492	5.05	0.000	.3604098	.8176661
_cons	-1.379241	.2688518	-5.13	0.000	-1.906181	-.8523008

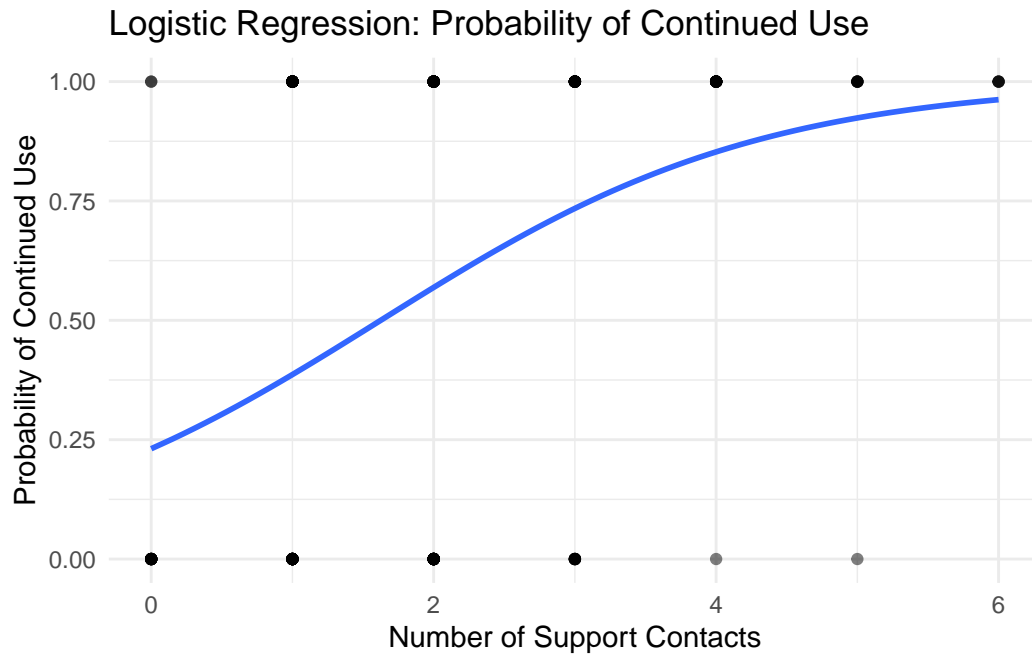
7.9.1 Output description

Term	Description
Coef.	Estimates of the regression coefficients.
Std. Err.	Standard errors of the coefficients.
z	The test statistic for each coefficient.
P> z	The p-value associated with each coefficient, indicating whether it is statistically significant. If the p-value is less than the significance level (typically 0.05), we reject the null hypothesis that the coefficient is equal to zero.

7.10 Plotting the Results in R

```
# Plotting the logistic regression curve
ggplot(data, aes(x = support_contact, y = continued_use)) +
  geom_point(alpha = 0.5) +
  geom_smooth(method = "glm", method.args = list(family = "binomial"), se = FALSE) +
  labs(title = "Logistic Regression: Probability of Continued Use",
       x = "Number of Support Contacts",
       y = "Probability of Continued Use") +
  theme_minimal()
```

`geom_smooth()` using formula = 'y ~ x'



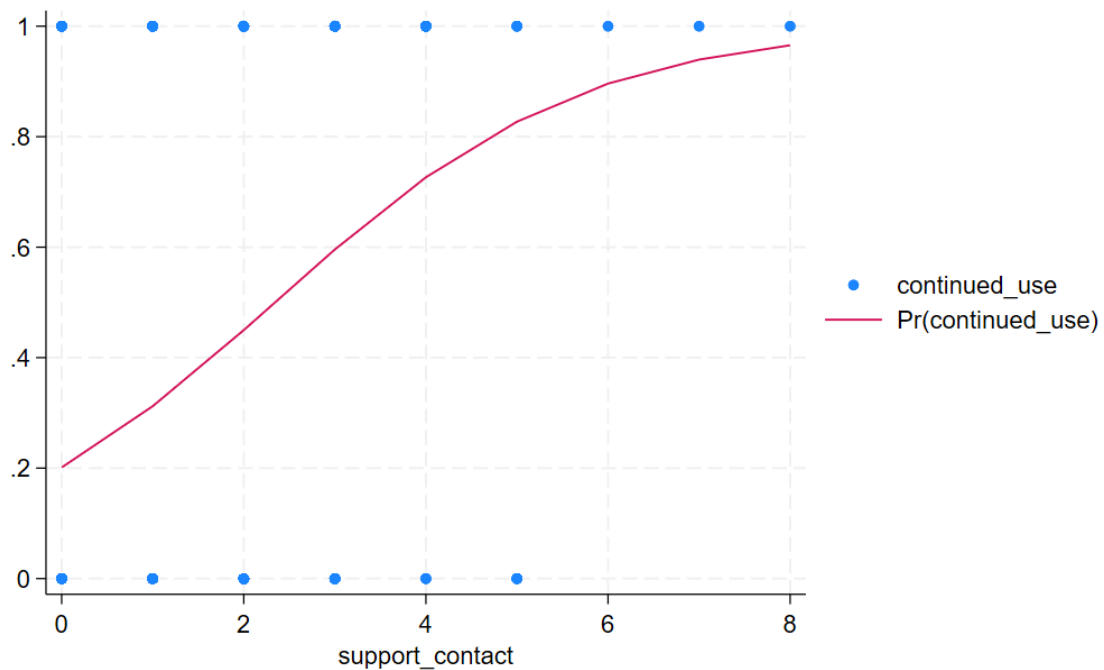
7.11 Plotting the Results in Stata

```
* Create logistic regression plot

/* Generate the predicted probabilities
This command generates the predicted probabilities from the logistic regression model and stores them in a new variable 'prob'
predict prob, pr

/* Sort the data by the predictor variable. Sorting the data by 'support_contact' ensures that the data is ordered by the number of support contacts
sort support_contact

/* Plot the scatter plot with the logistic regression line This command creates a scatter plot of 'continued_use' against 'support_contact' and overlays the predicted probabilities from the logistic regression model
twoway (scatter continued_use support_contact) (line prob support_contact)
```



7.12 Assumptions

Assumption	Description
Binary Outcome	The dependent variable should be binary.
Independence	Observations should be independent of each other.
Linearity of logit	The logit (log-odds) of the outcome should be linearly related to the predictors.
No multicollinearity	The predictors should not be highly correlated with each other.
Large sample size	Logistic regression typically requires a large sample size to provide reliable estimates.

7.13 R

7.13.1 Binary outcome

```
table(data$continued_use)
```



```
0    1
89 111
```

7.13.2 Independence

Verify that observations are independent. This is usually ensured by the study design.

7.13.2.1 Linearity of logit

```
library(car)
```

Loading required package: carData

Attaching package: 'car'

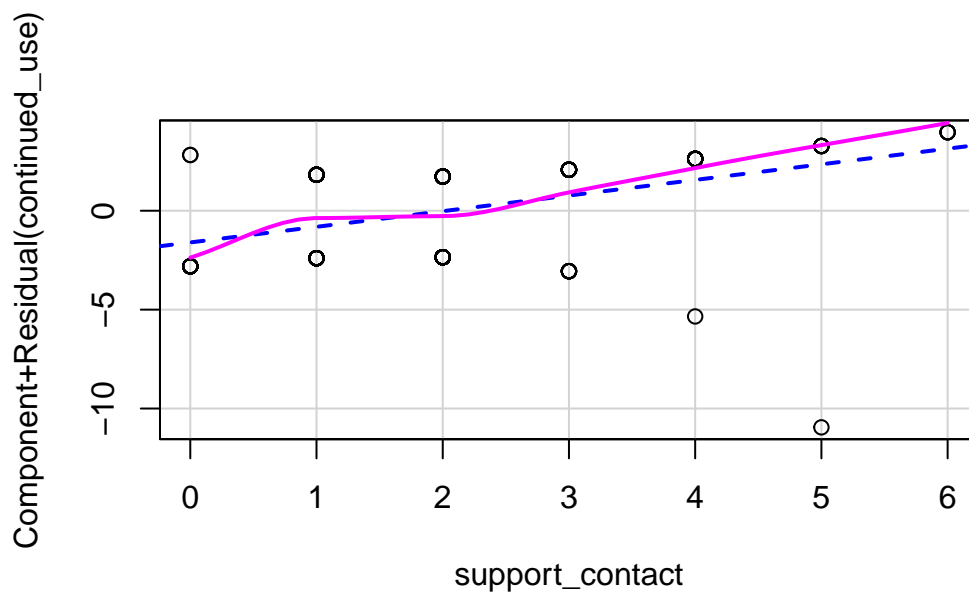
The following object is masked from 'package:dplyr':

```
recode
```

The following object is masked from 'package:purrr':

```
some
```

```
logit_model <- glm(continued_use ~ support_contact, data = data, family = binomial)
crPlots(logit_model)
```



7.13.3 No Multicollinearity

Note: The assumption of no multicollinearity is relevant only when you have at least two predictors in your model. Multicollinearity occurs when two or more predictors are highly correlated with each other, which can make it difficult to determine the individual effect of each predictor on the dependent variable.

Check for multicollinearity among predictors.

```
vif(YOUR_MODEL_HER)
```

```
Error in eval(expr, envir, enclos): object 'YOUR_MODEL_HER' not found
```

7.13.4 Large sample size

Ensure you have a sufficiently large sample size. A rule of thumb is at least 10 events per predictor variable.

7.14 Stata

7.14.1 Binary outcome

```
tabulate continued_use
```

```
no variables defined  
r(111);
```

```
r(111);
```

7.14.2 Independence

Verify that observations are independent. This is usually ensured by the study design.

7.14.3 Linearity of logit

```
gen logit_support_contact = log(support_contact / (1 - support_contact))  
scatter logit_support_contact support_contact
```

```
support_contact not found  
r(111);
```

```
r(111);
```

7.14.4 No Multicollinearity

```
estat vif
```

```
last estimates not found  
r(301);
```

```
r(301);
```

7.14.5 Large sample size

Ensure you have a sufficiently large sample size. A rule of thumb is at least 10 events per predictor variable.

7.15 Syntax Comparison: R vs Stata

This table summarizes the main differences between R and Stata in terms of syntax for performing Logistic Regression analysis.

Task	R Command	Stata Command
Simulating Data	<code>rpois()</code> , <code>rbinom()</code>	<code>rpoisson()</code> , <code>rbinomial()</code>
Setting Seed for Reproducibility	<code>set.seed(123)</code>	<code>set seed 123</code>
Visualizing Descriptives	<code>ggplot()</code> with <code>geom_boxplot()</code>	<code>graph box</code>
Running Logistic Regression	<code>glm()</code> with <code>family = "binomial"</code>	<code>logit</code>
Plotting the Results	<code>ggplot()</code> with <code>geom_smooth(method = "glm", ...)</code>	<code>twoway scatter</code> and <code>lfit</code>

8 Summary

In summary, this book has no content whatsoever.

1 + 1

[1] 2

References

Knuth, Donald E. 1984. “Literate Programming.” *Comput. J.* 27 (2): 97–111. <https://doi.org/10.1093/comjnl/27.2.97>.