# Experimental and Quasi-Experimental Designs

## for Generalized Causal Inference

Shadish | Cook | Campbell

# Contents

## 9. PRACTICAL PROBLEMS 1: ETHICS, PARTICIPANT RECRUITMENT, AND RANDOM ASSIGNMENT

## 13. GENERALIZED CAUSAL INFERENCE: METHODS FOR MULTIPLE STUDIES

## 14. A CRITICAL ASSESSMENT OF OUR ASSUMPTIONS

# Designing Experiments and Analyzing Data

## Second Edition

## A Model Comparison Perspective



Scott E. Maxwell
Harold D. Delaney

# Contents

## III  MODEL COMPARISONS FOR DESIGNS INVOLVING WITHIN-SUBJECTS FACTORS

## 11  One-Way Within-Subjects Designs: Univariate Approach  525

# IV  ALTERNATIVE ANALYSIS STRATEGIES

# 15  An Introduction to Multilevel Models for Within-Subjects Designs

## Appendixes

Gareth James
Daniela Witten
Trevor Hastie
Robert Tibshirani

# An Introduction to Statistical Learning

with Applications in R

Springer

# Contents

# Data Analysis Using Regression and Multilevel/Hierarchical Models

ANDREW GELMAN

JENNIFER HILL

CAMBRIDGE

# Contents

# An Introduction to R
## Second Edition

W. N. Venables, D. M. Smith and the
R Development Core Team

# Table of Contents

# Phil Spector

# Data Manipulation with R

# Contents

# Using R for Introductory Statistics

## Second Edition

John Verzani

# Preface

These notes are an introduction to using the statistical software package `R` for an introductory statistics course. They are meant to accompany an introductory statistics book such as Kitchens *"Exploring Statistics"*. The goals are not to show all the features of `R`, or to replace a standard textbook, but rather to be used with a textbook to illustrate the features of `R` that can be learned in a one-semester, introductory statistics course.

These notes were written to take advantage of `R` version 1.5.0 or later. For pedagogical reasons the equals sign, `=`, is used as an assignment operator and not the traditional arrow combination `<-`. This was added to `R` in version 1.4.0. If only an older version is available the reader will have to make the minor adjustment.

There are several references to data and functions in this text that need to be installed prior to their use. To install the data is easy, but the instructions vary depending on your system. For Windows users, you need to download the "zip" file , and then install from the "packages" menu. In UNIX, one uses the command `R CMD INSTALL packagename.tar.gz`. Some of the datasets are borrowed from other authors notably Kitchens. Credit is given in the help files for the datasets. This material is available as an `R` package from:

> `http://www.math.csi.cuny.edu/Statistics/R/simpleR/Simple_0.4.zip` for Windows users.
> `http://www.math.csi.cuny.edu/Statistics/R/simpleR/Simple_0.4.tar.gz` for UNIX users.

If necessary, the file can sent in an email. As well, the individual data sets can be found online in the directory

> `http://www.math.csi.cuny.edu/Statistics/R/simpleR/Simple`.

This is version 0.4 of these notes and were last generated on August 22, 2002. Before printing these notes, you should check for the most recent version available from

> the CSI Math department (`http://www.math.csi.cuny.edu/Statistics/R/simpleR`).

# Contents

Peter Dalgaard

# Introductory Statistics with R

Second Edition

# Contents

# A Handbook of
# Statistical
# Analyses
## Using R

# SECOND
# EDITION

**Brian S. Everitt** and **Torsten Hothorn**

# Contents

*Practical Recipes for Visualizing Data*

# R Graphics
# Cookbook

*Winston Chang*

# Table of Contents

# R Graphics

## Paul Murrell

# Contents

## II  GRID GRAPHICS

# List of Figures

# List of Tables