

my-ebook

Parmeshvar

2025-07-06

Table of contents

1 Introduction

Introduction

DR.Harsh Pradhan, [Phone: +91-9930034241 , Email: harsh.231284@gmail.com], [Institute of Management Studies, Banaras Hindu University](#), Address: 18-GF, Jaipuria Enclave, Kaushambhi, Ghaziabad, India, 2010

Interest: [Goal Orientation](#) [Job Performance](#) [Consumer Behavior](#) [Behavioral Finance](#) [Bibiliometric Analysis](#) [Options as Derivatives](#) [Statistics](#) [Indian Knowledge System](#),

[Orcid ID](#)

[Google Scholar](#)

[Youtube ID](#)

[Academic Profile](#)

Courses offered:

1. Free online course, four weeks (MOOC), enrollments open: Introduction to Bayesian Data Analysis
2. Short (four-hour) tutorial on Bayesian statistics, taught at EMLAR 2022: [here](#)
3. Introduction to (frequentist) statistics
4. Introduction to Bayesian data analysis for cognitive science
5. BDA cover

1.1 Lecture notes

Download from [here](#).

1.2 Moodle website

All communications with students in Potsdam will be done through [this website](#). # Schedule

Week	Lecture	Main Topic	Subtopic	Video	PDF Resource
Week 1	2	Descriptive Statistics	Central Tendency	Video	Week 2.pdf
	2	Descriptive Statistics	Measure of Variability	Video	Same as above
	3	Descriptive Statistics	Describing Data	Video	Same as above
	4	Descriptive Statistics	Probability	Video	Same as above
	5	Descriptive Statistics	Distribution	Video	Same as above
Week 3	1	Descriptive Statistics	Z Table (Normal Distribution)	Video	Week 3.pdf
	2	Descriptive Statistics	Measuring Divergence	Video	Same as above
	3	Inferential Statistics	Sample and Population	Video	Same as above
	4	Inferential Statistics	Model Fit	Video	Same as above
	5	Inferential Statistics	Hypothesis and Error	Video	Same as above
Week 4	1	Terms of Statistics	Terms of Statistics	Video	Week 4.pdf
	2	Terms of Statistics	T-Test	Video	Same as above
	3	Terms of Statistics	T-Test in Detail	Video	Same as above
	4	ANOVA	ANOVA	Video	Same as above
Week 5	1	ANOVA	Example of ANOVA	Video	Week 5.pdf
	2	ANOVA	Types of ANOVA	Video	Same as above

Week	Lecture	Main Topic	Subtopic	Video	PDF Resource
Week 6	3	Correlation	Introduction to Correlation	Video	Same as above
	4	Correlation	Regression (Part 1)	Video	Same as above
	5	Correlation	Regression (Part 2)	Video	Same as above
	1	Correlation	R Script for Regression	Video	Week 6.pdf
	2	Chi Square	Chi Square	Video	Same as above
	3	Chi Square	Chi Square Test	Video	Same as above
Week 7	4	Logistic Function	Regression Function	Video	Same as above
	5	Logistic Function	Distribution	Video	Same as above
	1	Time Series	Intro to Time Series	Video	Week 7.pdf
	2	Time Series	Conditional Probability	Video	Same as above
	3	Time Series	Additional Concepts	Video	Same as above
	4	Time Series	Distribution	Video	Same as above
	5	Time Series	Poisson Distribution	Video	Same as above
	6	Index Numbers	Price & Quantity Index	Video	Same as above
	7	Decision Environments	Risk/Uncertainty, Bayes, Trees	Video	Same as above
	8	Time Series Analysis	Components, Trend, Seasonality	Video	Same as above
	9	Time Series Analysis	Least Squares Method	Video	Same as above
	1	Effect Size & Documentation	Package/Library	Video	Week 8.pdf

Week	Lecture	Main Topic	Subtopic	Video	PDF Resource
2		Effect Size & Documentation	RStudio vs RKward	Video	Same as above
3		Effect Size & Documentation	Flexplot	Video	Same as above
4		Effect Size & Documentation	Functions	Video	Same as above
5		Effect Size & Documentation	R Shiny & R Markdown	Video	Same as above
6		Effect Size & Documentation	Application with Real Datasets	Video	Same as above
7		Effect Size & Interpretation	Importance in Testing	Video	Same as above
8		Effect Size & Interpretation	Installing dplyr, ggplot2	Video	Same as above
9		Effect Size & Interpretation	Visual Model Interpretation	Video	Same as above
10		Effect Size & Interpretation	Creating/Using Functions	Video	Same as above
11		Effect Size & Interpretation	Report, Dashboard, Interactivity	Video	Same as above

2 week 6

Table of Contents

Introduction__

Chi-Square Test of Goodness of Fit

Chi-Square Test of Independence

Non-Parametric Tests

Non-Linear and Logistic Regression

Poisson & Negative Binomial Distribution

Robust and Bayesian Regression

Model Fit Diagnostics

Exercises, Simulations, & Datasets

Summary

References

1. Introduction

This Week 6 eBook focuses on advanced statistical procedures for analyzing categorical and non-normal data using RKWard, a GUI-based frontend to R.

We address: - When traditional parametric methods fail - Tools for ordinal, non-linear, or count data - How to interpret diagnostic plots, residuals, and goodness-of-fit metrics

2. Chi-Square Test of Goodness of Fit

Theory Refresher

Use this test to see if observed frequency data matches a theoretical distribution (e.g., uniform, binomial, Poisson).

Example 1: Dice Fairness

```
obs <- c(9, 7, 6, 4, 5, 5) expected <- rep(sum(obs)/6, 6) chisq.test(obs, p = rep(1/6, 6))
```

 Example 2: Simulated Biased Die (Monte Carlo)

```
set.seed(42) sim_data <- sample(1:6, size = 600, replace = TRUE, prob = c(0.1, 0.1, 0.2, 0.2, 0.2, 0.2)) table_sim <- table(sim_data) chisq.test(table_sim, p = rep(1/6, 6))
```

 Example 3: Poisson-GOF for Counts

```
library(MASS) data_counts <- rpois(100, lambda = 3) obs_table <- table(data_counts)
exp_probs <- dpois(as.numeric(names(obs_table)), lambda = 3) chisq.test(obs_table, p =
exp_probs/sum(exp_probs)) Visualizing Frequencies
```

```
barplot(rbind(obs, expected), beside = TRUE, col = c("skyblue", "orange"), legend.text =
c("Observed", "Expected"), main = "Dice Roll Distribution")
```

3. Chi-Square Test of Independence
Purpose Test whether two categorical variables are independent.

Example 1: Gender vs Preference

```
df <- data.frame( Gender = c("Male", "Male", "Female", "Female"), Laptop = c("Gaming", "Non-
Gaming", "Gaming", "Non-Gaming"), Freq = c(27, 8, 5, 7) ) table_df <- xtabs(Freq ~ Gender +
Laptop, data = df) chisq.test(table_df)
```

Example 2: Titanic Survival

```
library(datasets) data(Titanic) chisq.test(Titanic) Example 3: Simulated Survey
set.seed(123) survey <- data.frame( Smoke = sample(c("Yes", "No"), 100, replace =
TRUE), Exer = sample(c("None", "Some", "Regular"), 100, replace = TRUE) ) tb <-
table(survey.Smoke, survey.Exer) chisq.test(tb)
```

Association Strength
library(vcd) assocstats(tb)
4. Non-Parametric Tests Why Use Them? Parametric assumptions (normality, equal variance) are not always met. Non-parametric tests allow analysis without these constraints.

Common Tests Parametric Non-Parametric Equivalent One-sample t-test Wilcoxon Signed-Rank Test Two-sample t-test Mann-Whitney U Test One-Way ANOVA Kruskal-Wallis Test Two-Way ANOVA Friedman Test Pearson Correlation Spearman Rank Correlation

Example 1: Wilcoxon Test (Single Sample)

```
data <- c(3.1, 3.6, 3.8, 4.0, 3.5) wilcox.test(data, mu = 3.5)
```

Example 2: Mann-Whitney (Between Groups)

```
group_a <- c(10, 12, 14, 16) group_b <- c(8, 9, 10, 11) wilcox.test(group_a, group_b)
```

Example 3: Kruskal-Wallis on Iris

```
kruskal.test(Sepal.Length ~ Species, data = iris)
```

Example 4: Spearman Rank Correlation

```
cor.test(iris.Sepal.Length, iris.Petal.Length, method = "spearman")
```

Next: Part 2 — covering:

Non-Linear Regression

Logistic Regression

Poisson & Negative Binomial

Robust & Bayesian Regression

Model Fit Diagnostics

Simulations, Interactive Plots

5. Non-Linear and Logistic Regression

5.1 Non-Linear Regression

Used when data shows curvature, not a straight-line relationship.

Example 1: Quadratic Fit

```
“r x <- 1:10 y <- 5 + 2 * x^2 + rnorm(10, 0, 10) model_quad <- lm(y ~ poly(x, 2, raw = TRUE)) summary(model_quad) plot(x, y) lines(x, predict(model_quad), col = “red”) Example 2: Exponential Growth
```

```
x <- 1:20 y <- 2 * exp(0.3 * x) + rnorm(20, 0, 10) df <- data.frame(x, y) model_exp <- nls(y ~ a * exp(b * x), data = df, start = list(a = 1, b = 0.1)) summary(model_exp)
```

Example: Student Pass/Fail

```
students <- data.frame( Hours = c(1,2,3,4,5,6,7,8,9,10), Pass = c(0,0,0,1,1,1,1,1,1,1) )
```

```
log_model <- glm(Pass ~ Hours, data = students, family = binomial()) summary(log_model) Predict Probabilities
```

```
studentsprob <- predict(log_model, type = “response”) plot(studentsHours, students$prob, type = “b”, col = “blue”) ROC Curve
```

```
library(pROC) roc_obj <- roc(studentsPass, studentsprob) plot(roc_obj) auc(roc_obj)
```

{r}

Test Fit

```
observed <- table(data_pois) expected <- dpois(as.numeric(names(observed)), lambda = lambda) chisq.test(observed, p = expected / sum(expected))
```

```
library(MASS) nb_data <- rnbinom(100, size = 5, mu = 4) hist(nb_data, col = “darkred”, main = “Negative Binomial”) Compare Fit
```

```
mean(data_pois); var(data_pois) # Poisson: mean variance mean(nb_data); var(nb_data) # NB: var > mean
```

```
library(MASS) x <- 1:10 y <- 2*x + rnorm(10) y[10] <- 100 # Outlier
```

```
model_rlm <- rlm(y ~ x) summary(model_rlm) plot(x, y) abline(model_rlm, col = “red”) Bayesian Regression (brms)
```

```
library(brms) data <- data.frame(x = rnorm(100), y = rnorm(100)) model_brm <- brm(y ~ x, data = data, family = gaussian(), chains = 2, iter = 1000) summary(model_brm) plot(model_brm)
```

AIC(model_quad, log_model) BIC(model_quad, log_model) Residual Plots

```
par(mfrow=c(2,2)) plot(log_model) Durbin-Watson Test
```

```
library(car) durbinWatsonTest(log_model)
```

chisq.test(Titanic) Challenge 2: Spearman on mtcars

cor.test(mtcarsmpg, mtcarshp, method = “spearman”) Challenge 3: Logistic + Polynomial

```
mtcarsam <- as.factor(mtcarsam) log_mod <- glm(am ~ poly(mpg, 2), data = mtcars, family = binomial()) summary(log_mod)
```

 Challenge 4: Negative Binomial Fit

```
library(MASS) data <- rnegbin(100, theta = 2) fit_nb <- glm.nb(data ~ 1) summary(fit_nb)
```

 10. Summary This module brought together:

Chi-Square Tests for independence and fit

Non-parametric alternatives to parametric tests

Logistic Regression for classification

Poisson and NB distributions for count data

Robust and Bayesian inference for resistant modeling

Diagnostics to ensure model quality

References

Dr. Harsh Pradhan, BHU Lecture Notes R Core Team (2024). The R Project for Statistical Computing. MASS, brms, car, vcd, performance, tidyverse packages Text: Field, A. (2013). Discovering Statistics Using R

Next Steps

Coming in Part 3:

Multinomial and ordinal logistic regression

Zero-inflated Poisson (ZIP) and hurdle models

Bootstrapping and permutation tests

RMarkdown interactivity: sliders, code widgets

Custom diagnostic dashboards

Expanded regression use cases: finance, healthcare, social science

Brute-force simulations, grid search tuning, multiple datasets

Data cleaning + wrangling using dplyr, janitor, and tidymodels

12. Advanced Logistic Models

12.1 Multinomial Logistic Regression

Used when the outcome variable has more than two categories (e.g., “Low”, “Medium”, “High”).

```
library(nnet) data(iris) irisSize <- cut(iris$Sepal.Length, breaks=3, labels=c("Short", "Medium", "Long")) model_multi <- multinom(Size ~ Sepal.Width + Petal.Length, data=iris) summary(model_multi)
```

 12.2 Ordinal Logistic Regression For ordered categories.

```
library(MASS) housing <- data.frame( Sat = factor(sample(1:3, 100, replace = TRUE), labels = c("Low", "Med", "High")), Infl = sample(1:5, 100, replace = TRUE), Type = sample(c("Tower", "Apartment", "House"), 100, replace = TRUE) ) model_ord <- polr(Sat ~ Infl + Type, data = housing, Hess=TRUE) summary(model_ord)
```

 13. Zero-Inflated and Hurdle Models 13.1 Zero-Inflated Poisson (ZIP) Used when count data has excess zeros.