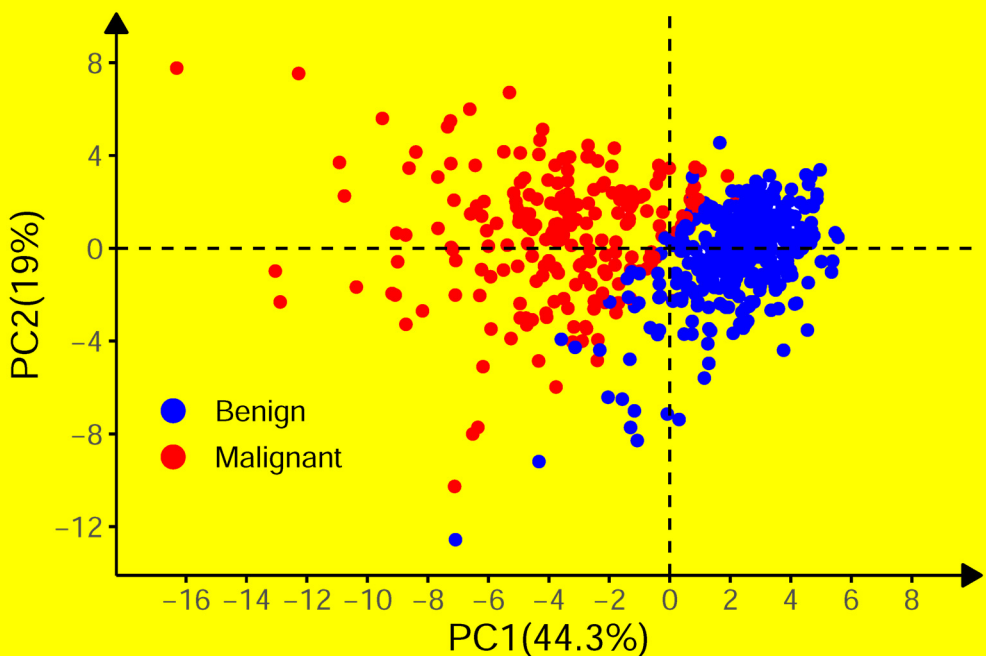


Introductory Principal Component Analysis Using R

A Practical Guide with RStudio



Kilem L. Gwet, Ph.D.

Get the entire book for \$9.99 using the link: <https://sites.fastspring.com/agreestat/instant/pca4rstudio/>

INTRODUCTORY PRINCIPAL COMPONENT ANALYSIS USING R

A Practical Guide with RStudio

Get the entire book for \$9.99 using the link: <https://sites.fastspring.com/agreestat/instant/pca4rstudio/>

Introductory Principal Component Analysis Using R

A Practical Guide with RStudio

Kilem L. Gwet, Ph.D.

AgreeStat Analytics
Gaithersburg, MD 20886-2696, USA

Copyright © 2023 by Kilem Li Gwet, Ph.D. All rights reserved.

Published by AgreeStat Analytics; in the United States of America.

No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or by an information storage and retrieval system – except by a reviewer who may quote brief passages in a review to be printed in a magazine or a newspaper – without permission in writing from the publisher. For information, please contact AgreeStat Analytics, at the following address:

AgreeStat Analytics
20416 Davencroft Ct,
Gaithersburg, MD 20886-4379
e-mail: gwet@agreestat.com

This publication is designed to provide accurate and authoritative information in regard of the subject matter covered. However, it is sold with the understanding that the publisher assumes no responsibility for errors, inaccuracies or omissions. The publisher is not engaged in rendering any professional services. A competent professional person should be sought for expert assistance.

Publisher's Cataloging in Publication Data:

Gwet, Kilem Li

Introductory Principal Component Analysis Using R

A Practical Guide with RStudio / By Kilem Li Gwet

p. cm.

Includes bibliographical references and index.

1. Statistics
2. Statistical Analysis
3. Data Analysis
4. Statistics - Study - Learning. I. Title.

Preface

Principal component analysis (PCA) will appeal to you if you have collected a large number of measurements on each of many subjects and find it difficult to extract useful information out of your dataset. Measurements are often associated with correlated variables, making it difficult to evaluate the impact of individual variables on subjects. PCA is a statistical technique that transforms the original correlated variables into a new set of *synthetic variables* also called *score variables* or *PC score variables* where PC stands for Principal Component. These variables have 2 interesting properties: (i) *they are uncorrelated*, and (ii) *they are ranked in such a way that the first few of them can explain most of the variation contained in the original dataset*.

You can see some key advantages of using only a few PC score variables as surrogates for the entire set of original variables:

- Since the first few PC score variables can explain most of the variation in the original dataset, your analysis can be based on one or two of them. The primary objective of PCA is dimensionality reduction. The 2 dominant PC score variables can give you an intuitive two-dimensional snapshot of a complex dataset.
- PC score variables are uncorrelated. Therefore, each can be interpreted independently from the others leading to a more stable statistical analysis.

The literature on principal component analysis (PCA) is abundant. This statistical technique is believed to have been introduced by **Pearson (1901)** before being rediscovered by **Hotelling (1933)**. While many computer scientists only discover this technique now in the era of machine learning, it is actually a very old statistical technique that has been used by statisticians for more than a century. Why this book?

This book can be used as supplement to many existing introductory books on PCA, which only cover specific aspects of this technique. For example, the book written by **Dunteman (1989)** provides a good balance between a mathematical treatment of PCA and a gentle introduction to key concepts. Some readers may still find it too mathematical, whereas others may find it too dated

to incorporate modern software solutions, which are essential these days for implementing PCA. The more advanced books such as that of Jolliffe (2002), tend to cover too much ground and to overemphasize the mathematical aspects of PCA.

The main objective of this book is two-fold:

- First, this book will show you how to perform PCA using the R software. Although this software comes with a documentation that shows you how to perform PCA, R outputs are not always clearly explained. Moreover, documentation files always assume that you have the appropriate background in PCA. However, the background information you obtain from textbooks will still not explain R outputs. This book will close that gap.
- Secondly, this book will provide the reader with a non-mathematical and yet rigorous presentation of the nature of Principal Component Analysis. The use of numerical examples is extensive and is the main approach adopted for introducing new concepts. You will learn as you do, or as you see me do.

After you complete the reading this book, I expect you to have an in-depth understanding of what principal components are, how to compute them using the R software, and how you can use them to improve the performance of other statistical techniques.

About the Author

I received my PhD in Mathematics from Carleton University's School of Mathematics and Statistics, Ottawa (Canada) in 1997. My specialization was the design and analysis of statistical surveys. In the past few years however, I wrote books and papers in the field of inter-rater reliability analysis (for some of my works in this field see https://www.researchgate.net/profile/Kilem_Gwet). After applying principal components to analyze multivariate inter-rater reliability data, I decided to write my first book on PCA. That book is more mathematical and uses MS Excel to implement the methods (see Gwet, 2020).

If you have comments or questions, please contact me at gwet@agreestat.com. I will respond to your inquiry as early as I possibly can.

Kilem Li Gwet, Ph.D.

Contents

Acknowledgment	ix
1 Basics of Principal Components	1
1.1 <i>Introduction</i>	2
1.2 <i>A Glimpse into Principal Component Analysis</i>	3
1.3 <i>Principal Components and Associated Scores</i>	6
1.4 <i>A More Realistic PCA Problem</i>	8
1.4.1 <i>The Principal Component Scores</i>	9
1.4.2 <i>Visualizing Key Principal Component Scores</i>	11
1.4.3 <i>Interpreting PC Loadings</i>	13
1.5 <i>Datasets</i>	14
1.6 <i>Concluding Remarks</i>	17
2 Overview of R and RStudio	19
2.1 <i>Introduction</i>	20
2.2 <i>RStudio Cloud</i>	20
2.3 <i>RStudio Desktop</i>	25
2.4 <i>Using RStudio</i>	28
2.4.1 <i>R Commands in RStudio</i>	28
2.4.2 <i>RStudio Workflow</i>	29
2.5 <i>Solving Problems with R</i>	32
2.5.1 <i>R Objects</i>	34
2.5.2 <i>Built-In Functions</i>	49
2.6 <i>Concluding Remarks</i>	53
3 Computing Principal Components with R	55
3.1 <i>Introduction</i>	56
3.2 <i>Calculating the Principal Components</i>	57
3.2.1 <i>Problem: Analysis of the MTCars Dataset</i>	58

3.2.2	<i>Exploring the PCA Results</i>	63
3.3	<i>Comparing <code>prcomp()</code> and <code>princomp()</code> Functions</i>	68
3.4	<i>Importance of Data Centering & Standardization</i>	73
3.5	<i>Concluding Remarks</i>	80
4	Visualization of Principal Components	83
4.1	<i>Introduction</i>	84
4.2	<i>Useful Packages</i>	87
4.3	<i>Scree & Cumulative Variance Plots</i>	89
4.4	<i>PCA Score Plot</i>	98
4.5	<i>Loadings and Variable Contribution Plots</i>	103
4.6	<i>Concluding Remarks</i>	112
5	Basic Use of Principal Components	115
5.1	<i>Introduction</i>	116
5.2	<i>Optimal Number of Principal Components</i>	116
5.3	<i>Predicting Scores for New Data</i>	122
5.4	<i>Reconstructing Original Data from Limited PCs</i>	129
5.4.1	<i>Anomaly Detection Using Principal Components</i>	129
5.4.2	<i>Data Reconstruction in Matrix Notations</i>	138
5.5	<i>Concluding Remarks</i>	140
6	Statistical Analysis Based on Principal Components	143
6.1	<i>Introduction</i>	144
6.2	<i>Identifying Relevant Variables with Principal Components</i>	145
6.3	<i>Identifying Outliers with Principal Components</i>	151
6.4	<i>Principal Component Regression (PCR)</i>	159
6.5	<i>Improving Classification with PCA</i>	188
6.6	<i>Concluding Remarks</i>	197
A	Performing Linear Discriminant Analysis in R	199
	Bibliography	207
	Author Index	209
	Subject Index	210

Acknowledgment

I want to express gratitude to my family for their support while working on this book. Most of all, I thank my wife Suzy and our three girls Mata, Lelna, and Addia. They have all graciously put up with my insatiable computer habits, my long workdays, and busy weekends over the past few years.

Kilem Li Gwet, Ph.D.
Maryland, USA: November 2023

Get the entire book for \$9.99 using the link: <https://sites.fastspring.com/agreestat/instant/pca4rstudio/>