# Statistical and Machine-Learning Data Mining

**Techniques for Better Predictive Modeling
and Analysis of Big Data**

Third Edition

# Statistical and Machine-Learning Data Mining

## Techniques for Better Predictive Modeling and Analysis of Big Data

### Third Edition

# Bruce Ratner

*This book is dedicated to:*

*My father, Isaac—always encouraging, my role model who taught
me by doing, not saying, other than to think positive.*◆

*My mother, Leah—always nurturing, my friend who taught me to love love and hate hate.*

*My daughter, Amanda—always for me, my crowning and most significant result.*

# *Contents*

# *Preface to Third Edition*

Predictive analytics of big data has maintained a steady presence in the four years since the publication of the second edition. My decision to write this third edition is not a result of the success (units) of the second edition but is due to the countless positive feedback (personal correspondence from the readership) I have received. And, importantly, I have the need to share my work on problems that do not have widely accepted, reliable, or known solutions. As in the previous editions, John Tukey's tenets, necessary to advance statistics, flexibility, practicality, innovation, and universality, are the touchstones of each chapter's new analytic and modeling methodology.

My main objectives in preparing the third edition are to:

1. Extend the content of the core material by including strategies and methods for problems, which I have observed on the top of *statistics on the table* [1] by reviewing predictive analytics conference proceedings and statistical modeling workshop outlines.

2. Reedit current chapters for improved writing and tighter endings.

3. Provide the statistical subroutines used in the proposed methods of analysis and modeling. I use Base SAS© and STAT/SAS. The subroutines are also available for downloading from my website: http://www.geniq.net/articles.html#section9. The code is easy to convert for users who prefer other languages.

I have added 13 new chapters that are inserted between the chapters of the second edition to yield the greatest flow of continuity of material. I outline the new chapters briefly here.

The first new chapter, Chapter 2, follows Chapter 1 (Introduction). The chapter is entitled *Science Dealing with Data: Statistics and Data Science*. If one were not looking, then it would appear that someone hit the delete key on statistics and statisticians and replaced them by science and data scientists. I investigate whether the recently minted term data science implies statistics is a subset of a more developed and expanded domain or if data science a buzzed-up cloaking of the current state of statistics.

Chapter 8, *Market Share Estimation: Data Mining for an Exceptional Case*, follows the Chapter 7 about principal component analysis (PCA). In this chapter, a market share estimation model, unique in that it does not fit the usual survey-based market share scenario, uses PCA as the foundation for estimating market share for a real exceptional case study. I provide the SAS subroutines used in building the market share model for the exceptional case study.

Chapter 11, *Predicting Share of Wallet without Survey Data*, follows the chapter on logistic regression. The everyday approach for predicting share of wallet (SOW) is with survey data. This approach is always met with reluctance because survey work is time-consuming, expensive, and yields unreliable data. I provide a two-step method for predicting SOW without data, by defining a quasi-SOW and using simulation for estimating total dollars spent. The second step uses fractional logistic regression to predict SOW_q. Fractional response regression cleverly uses ordinary logistic regression for dependent variables that

assume proportions or rates. I present a case study in detail as well as provide SAS subroutines that readers should find valuable for their toolkits.

Chapter 19, *Market Segmentation Based on Time-Series Data Using Latent Class Analysis,* follows the chapter on market segmentation via logistic regression. In this chapter, I propose the model-based clustering method of latent class analysis (LCA). The innovative strategy of this segmentation is in my use of time-series data. The times-series LCA model is radically distinctive and not equal, and it will prove to be a template for treating times-series data in cross-sectional datasets and the application of LCA instead of the popular data-based heuristic k-means. I provide SAS subroutines so that data miners can perform similar segmentations as presented, along with a unique way of incorporating time-series data in an otherwise cross-sectional dataset.

Chapter 20, *Market Segmentation: An Easy Way to Understand the Segments,* is well-placed after the LCA-based market segmentation. The literature is replete with clustering methodologies, of which any one can serve for conducting a market segmentation. In contrast, the literature is virtually sparse in the area of how to interpret the segmentation results. This chapter provides an easy way to understand the discovered customer segments. I illustrate the new method with an admittedly simple example that will not belie the power of the approach. I provide SAS subroutines for conducting the proposed technique so data miners can add this worthy statistical technique to their toolkits.

Chapter 21, *The Statistical Regression Model: An Easy Way to Understand the Model*, is an extension of the method of understanding a market segmentation presented in Chapter 20. Its purpose is to provide an easy way to understand the statistical regression model, that is, ordinary least squares and logistic regression (LR) models. I illustrate the proposed method with an LR model. The illustration brings out the power of the method, in that it imparts supplementary information, making up for a deficiency in the ever-relied-upon regression coefficient for understanding a statistical regression model. I provide the SAS subroutines, which serve as a valued addition to any bag of statistical methods.

Chapter 23, *Model Building with Big Complete and Incomplete Data*, follows the chapter that uses CHAD as a method for imputation. This chapter overhears missing data warning the statistician, "You can't win unless you learn how to accept me." Traditional data-based methods (complete case analysis), predating big data, are known to be problematic with virtually all datasets. These methods now open a greater concern as to their unknown ineffectiveness on big data. I propose a two-stage approach, in which modeling of response on complete-case data precedes modeling of response on incomplete-case data via PCA. The two models can be used separately or combined, depending on the goals of the task. I provide SAS subroutines for the proposed method, which should become a utile technique for the statistical model builder.

Chapter 24, *Art, Science, Numbers, and Poetry,* is a high-order blend of artwork, science, numbers, and poetry, all inspired by the Egyptian pyramids, da Vinci, and Einstein. Love it or hate it, this chapter makes you think.

Chapter 27, *Decile Analysis: Perspective and Performance*, complements the preceding chapter on assessment of marketing models. Marketers use decile analysis to assess predictive incremental gains of their response models over responses obtained by chance. I define two new metrics, response model decile-analysis precision and chance model decile precision, which allow marketers to make a more insightful assessment as to the incremental gain of a response model over the chance model. I provide the SAS subroutines for constructing the two new metrics and the proposed procedure, which will be a trusty tool for marketing statisticians.

Chapter 28, *Net T-C Lift Model: Assessing the Net Effects of Test and Control Campaigns*, extends the practice of assessing response models to the proper use of a control group (found in the literature under names such as the uplift or net lift model) instead of the chance model as discussed in Chapter 27. There is large literature, albeit confusing and conflicting, on the methodologies of net lift modeling. I propose another approach, the *net T-C lift* model, to moderate the incompatible literature on this topic by offering a simple, straightforward, reliable model that is easy to implement and understand. I provide the SAS subroutines for the Net T-C Lift Model to enable statisticians to conduct net lift modeling without purchasing proprietary software.

Chapter 34, *Opening the Dataset: A Twelve-Step Program for Dataholics*, has valuable content for statisticians as they embark on the first step of any journey with data. Set in prose, I provide a light reading on the expectant steps of what to do when cracking open the dataset. Enjoy. I provide SAS subroutines of the twelve-step program in case the reader wants to take a nip.

Chapter 43, *Text Mining: Primer, Illustration, and TXTDM Software*, has three objectives: First, to serve as a primer, readable, brief though detailed, about what text mining encompasses, and how to conduct basic text mining; second, to illustrate text mining with a small body of text, yet interesting in its content; and third, to make text mining available to interested readers, by providing my SAS subroutines, named *TXTDM*.

Chapter 44, *Some of My Favorite Statistical Subroutines*, includes specific subroutines referenced throughout the book and generic subroutines for some second-edition chapters for which I no longer have the data. Lastly, I provide some of my favorite statistical subroutines, helpful in almost all analyses.

If there are any corrections post-production of the text, I will post them to the errata link http://www.geniq.net/articles.html#section9.

## Reference

1. Stigler, S. M., *Statistics on the Table*, Harvard University Press, Cambridge, MA, 2002.

# *Preface of Second Edition*

This book is unique. It is the only book, to date, that distinguishes between statistical data mining and machine-learning data mining. I was an orthodox statistician until I resolved my struggles with the weaknesses of statistics within the big data setting of today. Now, as a reform statistician who is free of the statistical rigors of yesterday, with many degrees of freedom to exercise, I have composed by intellectual might the original and practical statistical data mining techniques in the first part of the book. The GenIQ Model, a machine-learning alternative to statistical regression, led to the creative and useful machine-learning data mining techniques in the remaining part of the book.

This book is a compilation of essays that offer detailed background, discussion, and illustration of specific methods for solving the most commonly experienced problems in predictive modeling and analysis of big data. The common theme among these essays is to address each methodology and assign its application to a specific type of problem. To better ground the reader, I spend considerable time discussing the basic methodologies of predictive modeling and analysis. While this type of overview has been attempted before, my approach offers a truly nitty-gritty, step-by-step approach that both tyros and experts in the field can enjoy playing with. The job of the data analyst is overwhelmingly to predict and explain the result of the target variable, such as RESPONSE or PROFIT. Within that task, the target variable is either a binary variable (RESPONSE is one such example) or a continuous variable (of which PROFIT is a good example). The scope of this book is purposely limited, with one exception, to dependency models, for which the target variable is often referred to as the "left-hand" side of an equation, and the variables that predict and/or explain the target variable is the "right-hand" side. This is in contrast to interdependency models that have no left- or right-hand side. I devote a chapter to one type of interdependency model, which is tied into a dependency model. Because interdependency models comprise a minimal proportion of the data analyst's workload, I humbly suggest that the focus of this book will prove utilitarian.

Therefore, these essays have been organized in the following fashion. Chapter 1 reveals the two most influential factors in my professional life: John W. Tukey and the personal computer (PC). The PC has changed everything in the world of statistics. The PC can effortlessly produce precise calculations and eliminate the computational burden associated with statistics. One need only provide the right questions. Unfortunately, the confluence of the PC and the world of statistics has turned generalists with minimal statistical backgrounds into quasi-statisticians and affords them a false sense of confidence.

In 1962, in his influential article, "The Future of Data Analysis" [1], John Tukey predicted a movement to unlock the rigidities that characterize statistics. It was not until the publication of *Exploratory Data Analysis* [2] in 1977 that Tukey led statistics away from the rigors that defined it into a new area, known as EDA (from the first initials of the title of his seminal work). At its core, EDA, known presently as data mining or formally as statistical data mining, is an unending effort of numerical, counting, and graphical detective work.

To provide a springboard to more esoteric methodologies, Chapter 2 covers the correlation coefficient. While reviewing the correlation coefficient, I bring to light several issues unfamiliar to many, as well as introduce two useful methods for variable assessment. Building on the concept of smooth scatterplot presented in Chapter 2, I introduce in Chapter 3 the smoother scatterplot based on CHAID (chi-squared automatic

xxviii

Preface of Second Edition

interaction detection). The new method has the potential of exposing a more reliable depiction of the unmasked relationship for paired-variable assessment than that of the smoothed scatterplot.

In Chapter 4, I show the importance of straight data for the simplicity and desirability it brings for good model building. In Chapter 5, I introduce the method of symmetrizing ranked data and add it to the paradigm of simplicity and desirability presented in Chapter 4.

Principal component analysis, the popular data reduction technique invented in 1901, is repositioned in Chapter 6 as a data mining method for many-variable assessment. In Chapter 7, I readdress the correlation coefficient. I discuss the effects the distributions of the two variables under consideration have on the correlation coefficient interval. Consequently, I provide a procedure for calculating an adjusted correlation coefficient.

In Chapter 8, I deal with logistic regression, a classification technique familiar to everyone, yet in this book, one that serves as the underlying rationale for a case study in building a response model for an investment product. In doing so, I introduce a variety of new data mining techniques. The continuous side of this target variable is covered in Chapter 9. On the heels of discussing the workhorses of statistical regression in Chapters 8 and 9, I resurface the scope of literature on the weaknesses of variable selection methods, and I enliven a notable solution for specifying a well-defined regression model in Chapter 10 anew. Chapter 11 focuses on the interpretation of the logistic regression model with the use of CHAID as a data mining tool. Chapter 12 refocuses on the regression coefficient and offers common misinterpretations of the coefficient that point to its weaknesses. Extending the concept of the coefficient, I introduce the average correlation coefficient in Chapter 13 to provide a quantitative criterion for assessing competing predictive models and the importance of the predictor variables.

In Chapter 14, I demonstrate how to increase the predictive power of a model beyond that provided by its variable components. This is accomplished by creating an interaction variable, which is the product of two or more component variables. To test the significance of the interaction variable, I make what I feel to be a compelling case for a rather unconventional use of CHAID. Creative use of well-known techniques is further carried out in Chapter 15, where I solve the problem of market segment classification modeling using not only logistic regression but also CHAID. In Chapter 16, CHAID is yet again utilized in a somewhat unconventional manner—as a method for filling in missing values in one's data. To bring an interesting real-life problem into the picture, I wrote Chapter 17 to describe profiling techniques for the marketer who wants a method for identifying his or her best customers. The benefits of the predictive profiling approach is demonstrated and expanded to a discussion of look-alike profiling.

I take a detour in Chapter 18 to discuss how marketers assess the accuracy of a model. Three concepts of model assessment are discussed: the traditional decile analysis, as well as two additional concepts, precision and separability. In Chapter 19, continuing in this mode, I point to the weaknesses in the way the decile analysis is used and offer a new approach known as the bootstrap for measuring the efficiency of marketing models.

The purpose of Chapter 20 is to introduce the principal features of a bootstrap validation method for the ever-popular logistic regression model. Chapter 21 offers a pair of graphics or visual displays that have value beyond the commonly used exploratory phase of analysis. In this chapter, I demonstrate the hitherto untapped potential for visual displays to describe the functionality of the final model once it has been implemented for prediction.

I close the statistical data mining part of the book with Chapter 22, in which I offer a data-mining alternative measure, the predictive contribution coefficient, to the standardized coefficient.

With the discussions just described behind us, we are ready to venture to new ground. In Chapter 1, I elaborated on the concept of machine-learning data mining and defined it as PC learning without the EDA/statistics component. In Chapter 23, I use a metrical modelogue, "To Fit or Not to Fit Data to a Model," to introduce the machine-learning method of GenIQ and its favorable data mining offshoots.

In Chapter 24, I maintain that the machine-learning paradigm, which lets the data define the model, is especially effective with big data. Consequently, I present an exemplar illustration of genetic logistic regression outperforming statistical logistic regression, whose paradigm, in contrast, is to fit the data to a predefined model. In Chapter 25, I introduce and illustrate brightly, perhaps, the quintessential data mining concept: data reuse. Data reuse is appending new variables, which are found when building a GenIQ Model, to the original dataset. The benefit of data reuse is apparent: The original dataset is enhanced with the addition of new, predictive-full GenIQ data-mined variables.

In Chapters 26–28, I address everyday statistics problems with solutions stemming from the data mining features of the GenIQ Model. In statistics, an outlier is an observation whose position falls outside the overall pattern of the data. Outliers are problematic: Statistical regression models are quite sensitive to outliers, which render an estimated regression model with questionable predictions. The common remedy for handling outliers is "determine and discard" them. In Chapter 26, I present an alternative method of moderating outliers instead of discarding them. In Chapter 27, I introduce a new solution to the old problem of overfitting. I illustrate how the GenIQ Model identifies a structural source (complexity) of overfitting, and subsequently instructs for deletion of the individuals who contribute to the complexity, from the dataset under consideration. Chapter 28 revisits the examples (the importance of straight data) discussed in Chapters 4 and 9, in which I posited the solutions without explanation as the material needed to understand the solution was not introduced at that point. At this point, the background required has been covered. Thus, for completeness, I detail the posited solutions in this chapter.

GenIQ is now presented in Chapter 29 as such a nonstatistical machine-learning model. Moreover, in Chapter 30, GenIQ serves as an effective method for finding the best possible subset of variables for a model. Because GenIQ has no coefficients—and coefficients furnish the key to prediction—Chapter 31 presents a method for calculating a quasi-regression coefficient, thereby providing a reliable, assumption-free alternative to the regression coefficient. Such an alternative provides a frame of reference for evaluating and using coefficient-free models, thus allowing the data analyst a comfort level for exploring new ideas, such as GenIQ.

## References

1. Tukey, J.W., The future of data analysis, *Annals of Mathematical Statistics*, 33, 1–67, 1962.
2. Tukey, J.W., *Exploratory Data Analysis*, Addison-Wesley, Reading, MA, 1977.

# *Acknowledgments*

This book, like all books—except the Bible—was written with the assistance of others. First and foremost, I acknowledge Hashem who has kept me alive, sustained me, and brought me to this season.

I am grateful to David Grubbs, my editor, who contacted me about outdoing myself by writing this book. I am indebted to the staff of the CRC Press/Taylor & Francis Group for their excellent work: Sherry Thomas, Editorial Assistant and Project Coordinator; Todd Perry, Project Editor; Victoria Jones, Copy Editor; Viswanath Prasanna, Senior Project Manager; Shanmuga Vadivu, Proofreader; Celia McCoy, Indexer, and Elise Weinger and Kevin Craig, Cover Designers.

# *Author*

**Bruce Ratner, PhD,** The Significant Statistician™, is president and founder of DM STAT-1 Consulting, the boutique firm for statistical modeling and analysis, data mining, and machine-learning. Bruce specializes in all standard statistical techniques, as well as recognized and innovative machine-learning methods, such as the patented GenIQ Model. Bruce achieves the clients' goals across varied industries: direct and database marketing, banking, insurance, finance, retail, telecommunications, healthcare, pharmaceutical, publication and circulation, mass and direct advertising, catalog marketing, e-commerce, Web-mining, B2B, risk management, and nonprofit fundraising.

Bruce's par excellence the expertise is apparent as he is the author of the best-selling book, *Statistical Modeling and Analysis for Database Marketing: Effective Techniques for Mining Big Data*. Bruce ensures the optimal solution methodology for his clients' marketing problems with rapid startup and timely delivery of project results. Bruce executes the highest level of statistical practice for his clients' projects. He is an often-invited speaker at public industry events, such as the SAS Data Mining Conference, and private seminars at the request of *Fortune* magazine's top 100 companies.

Bruce has his footprint in the predictive analytics community as a frequent speaker at industry conferences and as the instructor of the advanced statistics course sponsored by the Direct Marketing Association for more than a decade. He is the author of more than 100 peer-reviewed articles on statistical and machine-learning procedures and software tools. He is a coauthor of the popular textbook, *The New Direct Marketing*, and is on the editorial board of the *Journal of Database Marketing and Customer Strategy*.

Bruce is also active in the online data mining industry. He is a frequent contributor to *KDnuggets Publications,* the top resource for the data mining community. His articles on statistical and machine-learning methodologies draw a huge monthly following. Other online venues in which Bruce participates are the networking sites LinkedIn and ResearchGate, in which his postings on statistical and machine-learning procedures for big data have sparked countless rich discussions. Also, he is the author of his own *DM STAT-1 Newsletter* on the web.

Bruce holds a doctorate in mathematics and statistics, with a concentration in multivariate statistics and response model simulation. His research interests include developing hybrid modeling techniques, which combine traditional statistics and machine-learning methods. He holds a patent for a unique application in solving the two-group classification problem with genetic programming.

# 19

## *Market Segmentation Based on Time-Series Data Using Latent Class Analysis*

### 19.1 Introduction

Market segmentation is an often-used marketing model for efficient allocation of a company's resources. Market segmentation divides a population of customers into subpopulations—segments of customers. Customers within a segment are similar in their products and services, and customers across segments are dissimilar in their products and services. Market segmentation model implementation allows for effectively applying resources by targeting customers within their assigned segments. There are many statistical methods for market segmentation. A traditional and popular method is k-means clustering. A not-as-well-known method is latent class analysis (LCA). The purpose of this chapter is to present a novel approach to building a market segmentation model based on time-series data using LCA. I provide SAS© subroutines so data miners can perform segmentations similar to those presented, and I offer a unique way of incorporating time-series data in an otherwise cross-sectional dataset. The subroutines are also available for downloading from my website: http://www.geniq.net/articles.html#section9.

### 19.2 Background

I pedagogically outline this chapter. First, I concisely describe k-means clustering. Second, I cursorily review principal component analysis (PCA). I need to revisit PCA (from Chapter 7) because PCA and factor analysis (FA) are often confused, and FA is quite helpful in explaining LCA. Third, I present the LCA. Fourth, I compare LCA and k-means clustering. Lastly, I illustrate building a market segmentation model based on times-series data using LCA.

#### 19.2.1 K-Means Clustering*

K-means clustering generates k mutually exclusive groups by distances computed from one or more quantitative variables (Xs) [1]. Every observation belongs to one and only

---

* This section draws on https://support.sas.com/rnd/app/stat/procedures/fastclus.html.

one group of the k clusters. Because the number of clusters k is unknown, the model builder performs as many cluster solutions as desired. The typical k-means procedure uses Euclidean distance (i.e., least-squares calculations) that yields cluster means (of the observations in a cluster for the Xs). If groups exist, then all distances among observations in a group are less than all distances among observations in a different group.

The k-means algorithm is heuristic, which involves the following steps:

1. Diversely select k initial seeds (i.e., random points in the X-space of the observations).
2. These points represent initial cluster means.
3. Assign each observation to the cluster whose mean is the closest to the observation.
4. When all observations are assigned, recalculate the k means.
5. Repeat Steps #2 and #3, until the k means are stable.

## 19.2.2 PCA

PCA transforms a set of p variables* X1, X2, …, Xj, ..., Xp into p linear combination variables PC1, PC2, ..., PCj, ..., PCp (PCj denotes the j-th principal component). The essential objective of PCA is to establish a smaller set of the new PCj variables that represents most of the information (variation) in the original set of variables. An attractive analytic feature of the PCs is that they are uncorrelated with each other. The PCs are defined as:

$$PC1 = a11*X1 + a12*X2 +… + a1j*Xj + …+ a1p*Xp$$

$$PC2 = a21*X1 + a22*X2 +…+ a2j*Xj + …+ a2p*Xp$$

$$\vdots$$

$$PCi = ai1*X1 + ai2*X2 +…+ aij*Xj + …+ aip*Xp$$

$$\vdots$$

$$PCp = ap1*X1 + ap2*X2 +…+ apj*Xj + …+ app*Xp$$

where the aij's are the PC coefficients.

## 19.2.3 FA

There is confusion among too many statistics practitioners who do not understand the difference between PCA and FA. The reasons for the confusion is perhaps because:

1. PCA is often mentioned in textbooks as a particular case of FA.
2. Statistical computer packages treat PCA as an option in FA modules.

---

* For ease of presentation, the Xs are standardized.

3. PCA and FA both aim to reduce the dimensionality of the given dataset: PCA and FA are data reduction techniques. For example, a dataset with, say, 1,000 variables can be reduced to a statistically equivalent dataset with, say, only 150 variables.

4. PCA has been used extensively as a part of the FA solution.

### 19.2.3.1 FA Model

The FA model is defined as: p variables X1, X2, …, Xj, ..., Xp can be expressed (up to an error term) as a linear combination of m latent (unobservable) *continuous* variables or factor F1, F2, …, Fj, ..., Fm. Fs are defined as:

$$Xl = c11*Fl + c12*F2 + ... + c1j*Fj + ... + c1m*Fm$$

$$X2 = c21*Fl + c22*F2 + ... + c2j*Fj + ... + c2m*Fm$$

$$\vdots$$

$$Xi = ci1*F1 + ci2*F2 + ... + cij*Fj + ... + cim*Fm$$

$$\vdots$$

$$Xp = cp1*F1 + cp2*F2 + ... + cpj*Fj + ... + cpm*Fm$$

where cij's are like regression coefficients, called factor loadings.

The difference between PCA and FA is immediately apparent. Focus on the i-th PC and Xi:

$$PCi = ai1*X1 + ai2*X2 + … + aij2*Xij + … + aip*Xp$$

$$Xi = ci1*Fl + ci2*F2 + ... + cij*Fj + …+ cim*Fm$$

The PCs are a linear combination of the original *observed* Xs. Each Xi is a linear combination of *unobserved* latent factors F1, F2, …, Fj, ..., Fm.

FA attempts to achieve a reduction from p to m dimensions by postulating a model relating p Xs to m latent factors. In contrast, PCA directly transforms the Xs into PCs, which possess the desirable properties cited in Chapter 7.

### 19.2.3.2 FA Model Estimation

At first sight, the FA model (with simplified subscripts) looks like a standard ordinary least squares (OLS) regression model.

$$FA: X = c1*Fl + c2*F2 + ... + cj*Fj + ... + cm*Fm \tag{19.1}$$

$$OLS: Y = b1*X1 + b2*X2 + … bj*Kj + ... + bm*Xm \tag{19.2}$$

However, a closer inspection reveals a substantial difference. With FA, c's and Fs are both *unknown*. With OLS, b's are *unknown*, and the Xs are known.

The PCA solution to FA or the so-called principal FA involves the following steps:

1. Obtain initial estimates of the c's by taking the first m PCs from a PCA.
2. Obtain initial estimates of the Fs. X and c's are known. Thus, F can be initially determined.
3. Initial estimates are then fine-tuned, looping within Step #2, until optimal.

### 19.2.3.3 FA versus OLS Graphical Depiction

The typical graphical depictions of the FA and OLS models clearly indicate the theoretical framework of these two models in Figures 19.1 and 19.2.

The two models illustrated with three Xs are now visibly different in conceptual ways. The FA model has line arrows from the factor F to the Xs. The direction of the arrows indicates that the factor, the unobserved latent variable, affects the independent variables X1, X2, and X3. The lower case e's indicate the unknown errors associated with the measurement of the observed corresponding Xs. The OLS regression model has line arrows from the independent variables X1, X2, and X3 to the dependent variable Y. The direction of the arrows indicates the independent variables X1, X2, and X3 affect the dependent variable Y. There are no measurement errors assumed in OLS that influence Xs.

### 19.2.4 LCA versus FA Graphical Depiction

As a proper introduction to LCA is in the next section, I bring out the traditional setting of LCA, which is often considered similar to FA—with one difference. When the LCA model (Figure 19.3) is compared to the FA model (Figure 19.1), the similarity is visible. The significant and useful difference between the methods is the factor F of FA is a *continuous* latent variable, while the latent variable LC of LCA is *categorical*.

The usual graphic display for LCA is Figure 19.3. However, this display does not clearly portray the categorical nature of LC. I offer a refocused visual for an LCA 2-Class model in Figure 19.4, which clearly indicates LC is categorical. LC(1) and LC(2) are the two categories of LC. The crisscross-like arrows indicate each LC affects the independent variables X1, X2, and X3. Important to note, LCA nomenclature for an independent variable is *indicator*.



**FIGURE 19.1**
FA model graphic.

**FIGURE 19.2**
OLS model graphic.



**FIGURE 19.3**
LCA model graphic.



**FIGURE 19.4**
LCA 2-Class model graphic.

## 19.3  LCA

The concept of LCA, a statistical technique for identifying unobservable subpopulations within a population, was conceived in 1950 by Lazarsfeld, the father of this concept [2]. In 1968, Lazarsfeld (by then the *grandfather* of LCA) published the first comprehensive treatment of LCA with only categorical indicators, but he did not provide a reliable method for parameter estimation [3]. In 1974, Goodman solved the problem of obtaining maximum likelihood estimates[*] of LCA model parameters [4]. The traditional LCA was generalized by Vermunt in 1998 to include all scales of data—categorical, continuous, count, and ordinal [5].

### 19.3.1  LCA of Universal and Particular Study

As an illustration of how LCA works along with its output, I present a well-known and often revisited "Role Conflict and Personality" article, often referred to as the Universal and Particular Study, conducted in 1950 [6], which Goodman examined in 1974 [7]. In the 1950 study, 216 Harvard/Radcliffe undergraduates were asked how they would respond in four role-conflict situations. The premise of the role conflict study is "What right has your friend to expect you to protect him?" The four conflict scenarios[†] are:

1. You are riding in a car driven by a close friend, and he hits a pedestrian. You know that he is going at least 35 miles an hour (MPH) in a 20 MPH speed zone. There are no other witnesses. His lawyer says that if you testify under oath that the speed was only 20 MPH, it may save your friend from serious consequences.
   a. Universalistic response: "He has *no right* as a friend to expect me to testify to the lower figure."
   b. Particularistic response: "He has *a right* as a friend to expect me to give evidence to the lower figure."
2. As a physician, your friend asks you to "shade doubts" about a physical examination for an insurance policy.
3. As a drama critic, your friend asks you to "go easy on a review" of a bad play in which all of his savings are invested.
4. As a member of the board of directors, your friend intimates that you will "tip him off" about financially ruinous, though secret, company information.

   The response patterns of the 216 respondents generate the 16 (= $2 \times 2 \times 2 \times 2$) row patterns in Table 19.1.[‡] For each of the conflict scenarios (A, B, C, D), responses tending toward one that is universalistic are indicated by "+", and responses tending toward one that is particularistic are indicated by "–" [7].

#### 19.3.1.1  Discussion of LCA Output

Hagenaars and McCutcheon perform a simple two-class LCA on the Universal-Particular study [8]. The vital statistics of LCA output in Table 19.2 consist of (1) latent class probabilities and (2) conditional probabilities.

---

[*] Maximum likelihood estimation, specifically, the expectation–maximization (EM) process.
[†] The descriptions of the scenarios are taken from the original study.
[‡] Table is taken from Goodman 1974, p. 216.

**TABLE 19.1**

Response Patterns of the Four Conflict Scenarios

| A | B | C | D | Observed Frequency | A | B | C | D | Observed Frequency |
|---|---|---|---|---|---|---|---|---|---|
| + | + | + | + | 42 | − | + | + | + | 1 |
| + | + | + | − | 23 | − | + | + | − | 4 |
| + | + | − | + | 6 | − | + | − | + | 1 |
| + | + | − | − | 25 | − | + | − | − | 6 |
| + | − | + | + | 6 | − | − | + | + | 2 |
| + | − | + | − | 24 | − | − | + | − | 9 |
| + | − | − | + | 7 | − | − | − | + | 2 |
| + | − | − | − | 38 | − | − | − | − | 20 |

**TABLE 19.2**

LCA Output of Study

| Observed Indicators | Universal | Particular |
|---|---|---|
| Auto passenger friend | 0.993 | 0.714 |
| Insurance doctor friend | 0.939 | 0.329 |
| Drama critic friend | 0.929 | 0.354 |
| Board of directors friend | 0.769 | 0.132 |
| Latent Class Size | 0.280 | 0.720 |

The latent class probabilities are the sizes of the two latent classes, which are 0.280 and 0.720 for Universal and Particular, respectively. These probabilities show 28% and 72% of the study population are in the Universal and Particular classes, respectively.

Within the similarity of LCA and FA, a conditional probability is comparable to the factor loading in FA. Large conditional probabilities for a latent class imply the corresponding indicators are highly associated with the latent class and therefore define the latent class.

The large probabilities of the first class column, ranging from 0.993 to 0.769, confirm the correct labeling of the column as Universal. For the second class column, the large (0.714), moderate (0.329 and 0.354), and low (0.132) conditional probabilities reasonably define the Particular class. In other words, latent class Universal is reliably identified respondents, whereas Particular is moderately determined.

Within LCA, the conditional probability is the probability of an individual in a given class responding at a given level of the indicator. The conditional probability is 99.3% for a respondent in class Universal responding that there is no right to lie for the driver friend. The conditional probability is 71.4% for a respondent in Particular class responding that there is a right to lie for the driver friend.

Similarly, for a respondent who is in the Universal class and answers that there is no reason for doctor friend to lie, the conditional probability is 93.9%. And, the conditional probability is 32.9% for a respondent in Particular class responding that there is a reason for doctor friend to lie.

### 19.3.1.2 Discussion of Posterior Probability

The relevant statistic noticeably missing in Table 19.2 is the posterior probability—the probability of class membership *given* a respondent's answers to a given scenario. This statistic

allows for probabilistically classifying a respondent into one of the latent classes. The posterior probability is notationally represented as Prob (LC = c | A = i, B = j, C = k, D = l), where c = 1, 2; i = yes, no; j = yes, no; k = yes, no; and l = yes, no. The symbol | represents the word *given*.

The calculation of posterior probability needs the LCA joint probability, Prob (A = i, B = j, C = k, D = l, CL = c), defined as the product of:

- Prob (individual at CL = c)
- Prob (individual at A = i | CL = c)
- Prob (individual at B = j | CL = c)
- Prob (individual at C = k | CL = c)
- Prob (individual at D = l | CL = c)

Thus, the posterior probability of an individual belonging to latent class c (CL = c) given A = i, B = j, C = k, D = l is in Equation 19.3:

$$\mathrm{Prob}(LC = c \mid A = i, B = j, C = k, D = l)$$

$$= \mathrm{Prob}(A = i, B = j, C = k, D = l, CL = c) / \mathrm{sum}\left[\mathrm{Prob}(A = i, B = j, C = k, D = l, CL = c)\right], \text{across c}$$

(19.3)

Thus far, the discussion of LCA regards only categorical indicators. LCA with continuous and categorical indicators generates output virtually identical to that of LCA with only categorical indicators. Expectantly, the statistical computations of LCA with continuous variables add a level of detail and knowledge of the mathematical statistics underpinnings of LCA beyond the scope of this chapter. Accordingly, there is no discussion of the extended LCA, for which Vermunt (1998) is the definitive source.

## 19.4 LCA versus k-Means Clustering

LCA and k-means clustering have the same objective of dividing a population of individuals into k disjoint and exhaustive subpopulations (clusters) such that individuals within a group are as similar as possible, and the individuals among the groups are as dissimilar as possible. In statistics-speak, cluster construction seeks to maximize the between-cluster differences and minimize the within-cluster differences.

The fundamental underlying difference between LCA and k-means methodologies is that LCA is model-based, whereas k-means is a heuristic (technique). A heuristic is any approach to solving a problem that uses intuition based on the problem domain. Heuristics do not provide optimal solutions but rather good solutions. Markedly, LCA as a model-based statistical technique means that it posits a theoretical equation that represents the data-generating process of the population and drawn sample data. Moreover, the probability distribution of the statistical model is the distinguishing feature between model-based versus data-based heuristics.

K-means clustering is a popular, practical, and useful tool, especially for market segmentation. Similar to all techniques, k-means has strengths and weaknesses. The strengths (Items 1–4) and weaknesses (Items 5–12) are listed as follows:

1. K-means is easy to use, understand, and implement.
2. K-means performs best with many variables, which it can accommodate because it is computationally efficient.
3. K-means always produces a cluster solution.
4. K-means tends to produce tighter clusters than alternative techniques.
5. K-means can create a cluster with one or a handful of individuals.
6. K-means cannot suggest the optimal number of clusters.
7. K-means is sensitive to outliers. (When few individuals define a cluster, they are outliers.)
8. K-means is not a robust method in that pseudorandom seeds yield different solutions. Or, a holdout dataset can produce a different cluster solution.
9. The k-means algorithm does not reveal which variables are relevant.
10. The k-means algorithm is affected by variables with large variances. Accordingly, k-means often use standardized data.
11. K-means always produces a cluster solution. Always-a-solution gives a false sense of statistical security, in that k-means finds a good solution.
12. An objective criterion by which to assess the quality of the cluster solution does not exist.

LCA as a statistical model has the apparatus that indicates its strengths (first four items in the list that follows). Unfortunately, the apparatus does not eliminate weaknesses (last seven items in the list that follows).

1. There are objective criteria by which to assess the goodness-of-fit of the cluster solution.
2. Statistical testing exists for comparison between two or more candidate cluster solutions.
3. LCA solutions are not affected by the different scales and unequal or large variances of the indicators.
4. LCA, like all statistical models, produces residuals. Residual analysis is crucial in assessing remedies of lack-of-fit.
5. LCA cannot suggest the optimal number of clusters.
6. A serious problem with LCA is the assumption of conditional dependence, also known as local independence. Local independence means the indicators within a cluster class are independent of each other. The local independence is sometimes not a tenable assumption. The standard LCA model must be modified to account for this (http://john-uebersax.com/stat/faq.htm).
7. Methods for relaxing the conditional independence assumption are in-progress in recent years.
8. LCA uses maximum likelihood estimation like many statistical models. Thus, LCA solutions are subject to local maxima not necessarily the global maxima.

9. There are many goodness-of-fitness criteria, such as the basics: likelihood (L), log-likelihood (LL), and L-square (L²).

10. Some criteria weight the fit and the parsimony of a model based on sample size and degrees of freedom. They are:

    a. Akaike Information Criteria (AIC):* AIC, AIC(L²), AIC3(L²), AIC(LL), and AIC3(LL).

    b. Bayesian Information Criteria (BIC):† BIC, BIC2(L²), and BIC(LL).

    c. None of these criteria are universally deemed superior to another.

    d. The behavior of these statistics creates confusion and uncertainty as to which is the best cluster solution. For example, when comparing two models, it is not known how much difference in BIC is significant in order to choose one model over another [9].

11. Restrictions when assessing goodness-of-fit complicate building the model. Hypotheses are tested by imposing restrictions and determining how these limitations affect the fit of the model to the data. Two such restrictions include (1) equality constraint (e.g., parallel indicators or equal error rate) and (2) deterministic (e.g., setting conditional probability to a particular value—usually 1 or 0).

12. Sparseness, due to many indicators with many response options, leads to difficulties in model evaluation (e.g., determining the degrees of freedom).

## 19.5 LCA Market Segmentation Model Based on Time-Series Data

I present the building of a market segmentation model based on times-series data using LCA.

The proposed model is not to be confused with partitioning of a times series, producing a sequence of discrete segments to reveal the underlying structure of input time-series data. A typical method of a times-series segmentation is a piecewise linear regression, which forecasts, say, stock market trading, inputting time series into k straight lines, each of an equal length. The proposed LCA market segmentation model is a novel and efficient segmentation technique, based on time-series data that are used to create indicators, as required by LCA.

I outline the intended LCA time-series market segmentation model pedagogically in a step-by-step manner to show (1) the components of the time-series data preparation with the SAS subroutines used and (2) building the LCA segmentation. The LCA program used (not provided here) is among several commercial software packages available.

### 19.5.1 Objective

Hi-tech company PirSQ wants to build a market segmentation as an efficient way to target their best customers. The segmentation, providing segments or clusters, allows PirSQ to develop marketing strategies to optimize future unit orders.

---

* AIC= 2*Npar − 2*ln(L). AIC(L²) = L² − 2*df. AIC3(L²) = L² − 3*df. AIC(LL) = 2 log L + 2*Npar. AIC3(LL) = 2 log L + 3 Npar. Npar = number of parameters.

† BIC= -2*ln(L) + Npar*ln(N). BIC(L²) =L² − log(N)*df. BIC(LL) = 2*log L + log(N)*Npar. Npar = number of parameters, and N = sample size.