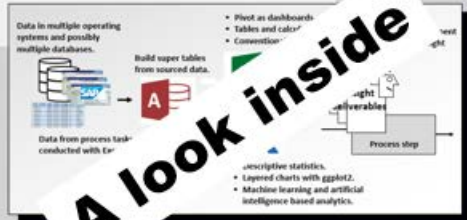**Data and Analytics Skills for Your Career Security**

Keeping it simple. . .
only the skills you're likely to use

Richard G. Lamb

Excerpt:
Table of Contents, Preface and Index

Additional "Look Inside" at https://analytics4strategy.com/book-look-inside
Book is available from Amazon
Paperback, 508 pages, 300 figures and outputs, 7.5 x 1.15 x 9.25 inches, 2.37 pounds

# Contents

# Preface

Does the steady drumbeat for data and analytics in your newsfeed make you uneasy? Are you worried that data and analytics will be the new Excel for which skills went from optional to mandatary for everyone's career? Are you searching for a way in but have not yet found a door other than some prohibitively huge and long undertaking to reach critical mass?

In other words, should you be worried for your career health? Yes, if your role engages with operating systems and Excel. And yes again, if your role in a managerial position entails decision-making upon the interpretation of standard reports from the operating system or Excel deliverables from other role holders.

You should be worried. As a few of these role holders progressively bring data and analytics skills into their roles, they will set a new standard that is immediate, visible and significant to the enterprise. Without the same skills you will not be able to meet the standard. Better yet, the best assurance of career security is to be one of those who sets the standard.

Maybe this is not new bad news to you. You see it coming. Your problem is that you have not yet stumbled across a way into the skills, short of awful. However, there is good news. Because I decided to take the long and awful path, this book is your serendipitous good luck. My bad luck was that a book such as this had never been written.

Why the book is your good luck is evident by my experience up to writing it. As an operational role holder, rather than data scientist, I had to learn much of everything there is to know about data and analytics before I could reduce the explanation of data and analytics to what you, in your operational role, need to know.

It took me six years to learn much of what there is to know. It then took me another more than a year to reduce what I had to learn to what you need to know to vaccinate your career against obsolescence and decline.

## Philosophy and Approach

The book is written to explain and demonstrate data and analytics as methods that are relevant to all enterprise operations rather than the universe of all human endeavors. The strategy works because all operations must deal with the same issues and sub-operations. Every operation has, in some form, requirements for setting direction, aggregate planning,

budget and variance control, planning, scheduling and executing the core task, staffing and optimizing the retention time of the core task in the stages through the operation.

It follows that all operations will thrive on the same data and analytics skills. Regardless of the type of operation, it is obviously impactful that its role holders be able to extract, explore, cleanse and mold to purpose the data captured in their operating systems. Regardless of the type of operation, it is obviously impactful that role holders be able to augment their experience and judgement with the insight that data and analytics make possible. Furthermore, the best augmenting insights are not and never will be available from the operating systems their roles depend upon.

Accordingly, a full explanation by demonstration of what you need to know, must necessarily be in the context of an operational domain. Therefore, as a context to all explanations by demonstration, the book has selected a domain in the management of manufacturing plants as a production asset. The type of operation is known as asset management. For those in the field of asset management, the book is a doubly serendipitous good fortune.

Chapter 1, Operational Availability is the Purpose, introduces asset management with respect to how it plays in the enterprise. It then establishes the measures of performance and the driving aspects to the measures. Because data and analytics should ride on purpose, you are advised to think out the equivalent framing of the operation in which you hold a role and the operations around it.

Of course, data and especially analytics are not, nor can be made, easy subjects. That is even after being bound within only what you need to know rather than all there is to know. As someone who read the book said, "You have to shut the door and turn off the TV."

However, you will not need to go far before you begin experiencing the payback of your concentration. At each chapter you will return to work with an actionable idea to upgrade the output of your role and the roles of others.

Let's talk about approaching the book to maximally bring about immediate payback to your organization and you as well. The book is written as the sequence of bringing the practices of asset management to be data driven. Asset management practitioners have got to love that. However, there is alternative path through the chapters for maximal "next-day" payback upon each stepwise accumulation of skills.

You should begin with Chapter 2, Data, Analytics and Software to be Data Driven. You need to know the territory. Just as important, you will discover that you already have at hand everything you need without first engaging management for permission to act on your payback idea or plead with them to buy you specialized software. Also important, the chapter will establish the definitions and demarcations with which you will be able calm

management down by explaining why it is that data drivenness entails learned skills with immediate operational ramifications rather than capital and costs for exotic software and endeavors.

Here is one thing I discovered while learning much of what there is to know about data and analytics. Most of the power of data-drivenness is achievable with only the skills of working with data. Only a small number of ideas, less than 20 percent, call for analytics. In other words, if the forces of the universe commanded that I was to be forever limited to what I could make happen with data, I would be disappointed but still excited about what I have been left with.

## Data Skills and Methods

When you know how to extract, explore and mold data into a dataset for purpose, your operation will immediately jump forward in its ramification for the enterprise. As a thought exercise, while you work through Chapter 3, Super Table from Operational Data, imagine what your operation would look like and perform like if every role holder were able to skillfully work with data as normal to their role. In fact, it is a readily achievable vision.

Working through Chapters 3, 8, 9 and 11, you will see in action almost every technique of building what the book calls super tables. The techniques are demonstrated in MS Access because the software is available to everyone by virtue of their organization's MS Office license. Furthermore, the principles of data that are demonstrated in Access travel easily to other software such as Tableau and Power BI.

Chapter 3, Super Tables from Operational Data, explains how data is extracted from operating systems, pulled into Access and joined together in a super table of all related variables. This is in contrast to Excel which leaves us divided and conquered because our data, although related, is present to us in disconnected tables. The chapter then explains how to create summary variables in the super table, a technique which loom huge in the ability to build insight from the data.

Chapter 8, Achieve Entirety of Data, is an important chapter because we must make all of our operational data available to building super tables. The chapter presents three frequently observed obstacles to entirety. To confront them, you will learn what changes must be made to operational roles and systems as necessary to achieve entirety.

Chapter 9, Workload-Based Budget and Variance, and Chapter 11, Budget-Based Schedule and Craft Capacity, demonstrate data skills in the context of expanding super tables to perform tedious, complex operational tasks. The value of data skills is highlighted by demonstrating a collection of related organizational tasks that are not otherwise possible in asset management and their absence has been a big impediment to enterprise

performance. It will jump out at you that complex laborious reporting and control tasks can be reduced to merely updating the source tables at the input side of the super table.

## Analytics Skills and Methods

What will be your career health at this juncture of acquired skills? It will be excellent. However, reaching preeminence would be even more excellent. More enticing is that preeminence is within your grasp because now "you know data."

Some time ago, I met a person in a large corporation whose meal ticket was his ability to produce Excel pivot charts and tables. The tables were not even super tables, nor did he know how to build them. Furthermore, the graphic capability of Excel to the fellow's claim to fame was based on single-perspective charting that was invented over 100 years ago, some as far back as the 1600's. Dashboards are an attempt to get past the limitation by placing multiple single-perspective charts in a single display. This is in contrast with layered perspectives in a single chart.

Thumb through Chapter 5, Layered Charting to Know Thy Data. You will see in action what is ggplot2 in the R software. The differences and ramifications will jump out at you. Imaging pulling such power to visualize data and measures into your role. In contrast, my buddy with his Excel pivots would be lost in your dust.

However, there is a bit of bad news. To get to layered visualization, you must step into the R software. As you will know from Chapter 2, R is a powerful and open software available for you or your firm to download without cost, limitations or strings.

The purpose of Chapter 4, The R Software in Action, is to get you up and running in R. Once again, the chapter is used to explain real life insights by demonstration. The chapter demonstrates the procedure to inspect a dataset statistically, reveal missing data, test variables for normal distribution and inspect the correlations between variables.

The philosophy of the book is that frequently occurring code will emerge from the explanations. Additional commonly occurring coding will emerge in the code for layered charting with the previously mentioned ggplot2 as well as from all of the other chapters to explain analytics. And just as for the Access code of the previously introduced chapters, you are left with R templates to substitute in your own variables, and R skills and code you will need to know.

You will arrive at another milestone upon working through Chapters 4 and 5. You will have accumulated the chops to begin moving into analytics. This is good because there is still an elephant in the room that requires analytics. It is to find and cleanse bad data. Of course, we should think of our best solutions for bad data are action taken to prevent a future of bad data in the operating system.

If there is bad data, what will be our strategy to rectify it? Chapter 6, Unearth and Rectify Bad Data, presents the types, decisions and schemes to rectify bad data. Several straightforward methods for some types of bad data were offered by Chapter 3. All others will require analytics to find and rectify them.

When some cases to a variable are bad, the question is what should they be? The cleansing process will lean on methods to replace bad cases, if necessary. However, the analytic is a core type for much more than replacing bad data with good. It is regression analytics to determine how strongly, if at all, specific operational variables are related to an outcome variable. Stated in the context of cleansing, the "outcome" variable can be the one with bad cases and we use the regression to estimate what they should be.

The chapters that explain two types of regressions are next in line to strengthen your career security. The analytics for relationships play in advancing operational effectiveness. For an outcome measure of performance, identifying the variables that matter most, or least, is powerful insight for assuring that your role is doing the right things right.

There are three mainstream regressions: linear, logistic and Poisson. The book will demonstrate linear and logistic. Chapter 7, Relate Operational Variables to Outcome, explains linear regression. Chapter 14, Recover Lost Classifications, explains logistic regression.

Regressions are a big step forward in your skills because they provide one of the most fundamental insights in operational capability. However, regression is not as simple as it would seem. We do not push all seemingly relevant variables into the model, run it and read the answer. Instead, there are considerable steps to choose and evaluate variables, and confirm fit: defined as how well the model accurately predicts the data.

Chapter 7 explains the process step by step. The same process generally applies to all regressions. I have never seen it mentioned in traditional texts that we can use the validation processes of regression to find outliers to the regression rather than only outliers to its variables. These are cases that the model would have never predicted and cases that have too much influence on the parameters of the model. Both are important insight to an operation because outliers may be telling us were to question and improve our operational processes and controls.

The difference between linear and logistic regression is their outcome variable. Linear predicts continuous numeric variables, whereas, logistic predicts classifications. Chapter 7 explains in detail the structure of the models and the distinction between them that rests upon the mathematics of prediction.

Chapter 14, Recover Lost Classifications, explains logistic regression while demonstrating a purpose for it beyond exploring relationships and outliers to outcomes. It

can be used to determine what an incorrect classification to a variable should have been. The strength of relationship of predictive variables to good classifications is used to cleanse the bad cases.

However, the purpose of the chapter is dual rather than merely to explain logistic regression. As the title suggests, the overarching purpose of the methods demonstrated is to recover lost classifications.

Accordingly, the chapter will also introduce and demonstrate the method known as naïve Bayes probability. We use naïve Bayes to determine from unstructured free-text variables, such as descriptions and notes, what a classification should be upon the probable occurrence of words in the text.

The next logical extension of your skills is to be able to prove there is a difference, or will be, as the consequence of change, improvement and enforcement of operational procedures and resources. Chapter 13, Prove There is a Difference, explains the body of analytics to make the determination vis-a-vis the operational situation. The analytics are a set of eight models from the types of two-means t test, ANOVA and multilevel.

Chapter 10, Through the Lens of Time Series, introduces another method to spot change through data and analytics skills. We look for patterns in variables presented in time series. However, our interest is not limited to change because overall we want to know what has happened over time.

Time series analytics are not to be confused with a line chart in Excel in which essential insights are lost in the simplicity. The chapter will explain by demonstration how to separate recurring cycles from the core pattern. It will explain how to measure the degree that one period reflects one or more previous periods. It will explain how the same principle is used to identify variables with a lead and lag relationship.

Just as importantly, we need to assure that a core pattern to a series is deterministic rather than a random walk. This is because a random walk can look deterministic. If not deterministic, we need to know to not devise actions or attempt forecasts as if it were.

Finally, in snaking through the chapters, we arrive at a final fundamental characteristic of any operation. It is the time it takes for each discrete core operational task to pass through the stages of the operation. Chapter 12, Elapsed Time Through Stages, explains by demonstration how to determine the statistics of retention in a stage and the chance of exit having remained in the stage for some duration. Just as important, the analytic enables us to determine which variables in the operation are most strongly related to retention and exit, thus, are the levers to reshape both characteristics of the stages that matter most.

## Where to Go Next

The book has bound its scope to the de facto favorite mainstream analytics and, for them, what you need to know to build and interpret them. References are provided if you want to take the watch apart rather than merely know how to tell time. What is most exciting is that, with each chapter, your ability to reach into, understand and engage additional methods outside the book's scope will grow and solidify.

I was once told that to qualify as an academic textbook there must be assignments with each chapter. We could say that replicating the demonstrations in each chapter is your assignment. My personal experience is that a powerful solidifying exercise is to emulate each demonstration and explore its code command by command. However, there is an assignment we could make for each chapter.

First, you must ask yourself a question. For which roles of my operation do the chapter's skills and methods apply, and how and why?" Second, you must act. The assignment is to return and install new methods in your role, thus, equip your organization with actionable insights it has never had before. Third, you must share the knowledge. The assignment is to pass each of your newly acquired skills to others by identifying roles for which the chapter's methods could be impactful and helping the holders of the roles to increase their career security by bringing new value to the enterprise through their work in the role.

Welcome to the new world that opens to you as a new age worker with all career security appertaining thereto.

## Data and Code to the Book

From literature and the internet, there are massive resources to explain all principles and methods using the R software and, for that matter, just about any software. The standing rule in the R community is that every provided explanation from any source must be accompanied with an example, code and dataset.

The book honors the rule by making the datasets and R code (scripts) to the chapters available to download from the webpage, https://analytics4strategy.com/data-and-code. You will need the datasets to emulate the demonstrations. The scripts are placed in the text for explanation by demonstration. However, since the code cannot be copied from the book and pasted into an R session, you have the option to copy and paste scripts into your R session.

## Color and Format

The book is a complex manuscript and, thus, presents challenges to formatting and cost. The body of the book is a progression through text explanations, supporting figures, output exhibits and the R and Access code that generated the exhibits. In formatting the pages, all demand a degree of sequential rigidity for being laced together in an explanation by demonstration that can be easily followed.

One issue in formatting is that some output exhibits are visualizations that depend heavily on color to communicate the insight they are coded to give us. This is especially the case for visualizations that demonstrate the advance graphics of the ggplot2 package of the R software.

Unfortunately, the cost to produce the book in color is prohibitive to pricing. To keep the price reasonable, the book is produced in grayscale. The readers, of course, can view each graphic in full color by running the provided R code. Alternatively, the reader can view any visualization in full color at webpage, https://analytics4strategy.com/book-look-inside.

Is it said in data and analytics, "the devil is in the details." When something does not work, it is often a tiny error in typing. Given the number of code blocks and internet addresses, hyphenation has not been used in formatting the text. Accordingly, readers need not wonder if a hyphen in a line of code or an internet address are from formatting or is as it should be. However, this policy occasionally creates some ugly lines in the text that hyphenation would eliminate. I ask for your forbearance.

In a few cases, the placement of output exhibits causes an over large empty white space at the bottom of a page. It is unavoidable given that some sequences must be honored with respect to the associated Access and R code upon which they are generated. Once again, I ask for your forbearance.

So, now onward we go through the fog and the dark where skill and education win over ignorance and superstition every time.

Richard G. Lamb, PE, CPA.

# Index

# Data and Analytics Skills for Your Career Security
## *Keeping it simple, only the skills you're likely to need*
### Richard G. Lamb

For those of us who are role holders in enterprise functioning, the personal purpose of acquiring practical working skills in data and analytics is to be able to better do what we already do and find new ways to do better yet. It follows that if you are a role holder who brings and incorporates data and analytics methods in your thoughts and tasks, your career outlook will be more secure and exciting. The book is written to be your gateway to the skills and to be the templates with which you will install the methods in your operational roles.

We all know that the field of data and analytics is huge and intimidating. It is a long slog to becoming comfortable. During the author's own long slog until arriving at the book, something exciting bubbled to the surface. There is a big difference between what we need to know and everything there is to know. We need to know what is possible as insight for decisions and functioning, we need to know how to get to the insight and, finally, we need to be able to interpret the insight. Just as the book does, we can leave the rest to the data scientists.

**About the Author:** In 2003, Richard Lamb, while struggling to get at the history captured in the databases of operational systems, found the skills to extract datasets of related history and join them in a super table of variables to make possible what was being envisioned for operational effectiveness. In 2014, Richard realized that, with statistical analytics and free enabling powerful pc-level software, an enterprise could ask and answer questions of operational effectiveness that are otherwise not possible. His activism to bring the epiphanies into the careers of role holders in the mainstream of operations has arrived at this book to explain data and analytics through the demonstration of methods.

Richard is a Registered Professional Engineer and Certified Public Accountant. He has previously authored two books: Availability Engineering and Management for Manufacturing Plant Performance, and Maintenance Reinvented for Business Performance. He has a BSCE, BBA and MBA from the University of Houston and a graduate-level Applied Statistics Certificate from the Texas A&M University.

https://analytics4strategy.com/data-and-code

**Analytics4Strategy**