

FIGURE 12-20 Frequency distribution with log of user score

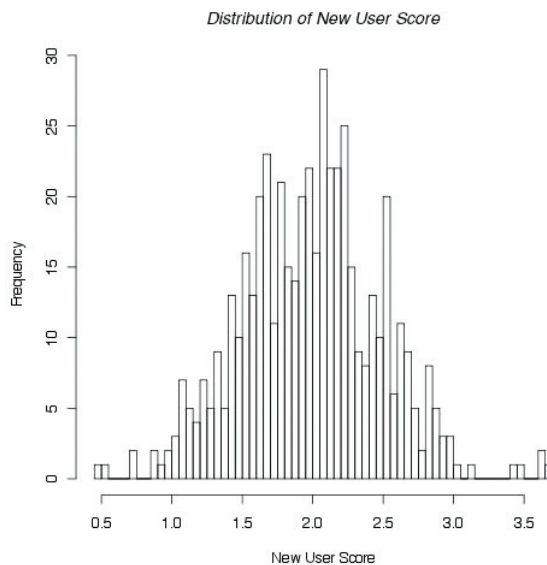


FIGURE 12-21 Frequency distribution of new user scores

Another idea may be to analyze the stability of price distributions over time to see if the prices offered to customers are stable or volatile. As shown in a graphic such as Figure 12-22, the prices appear to be stable. In this example, the user score of pricing remains within a tight band between two and three regardless of the time in days. In other words, the time in which a customer purchases a given product does not significantly influence the price she is willing to pay, as expressed by the user score, shown on the y-axis.

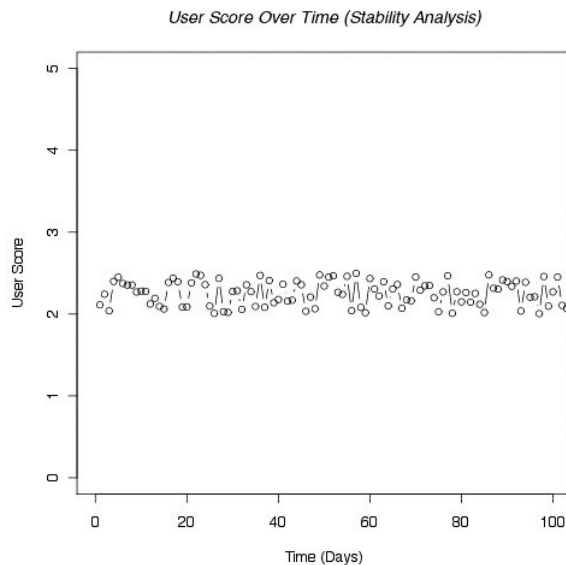


FIGURE 12-22 Graph of stability analysis for pricing

By this point the data scientist has learned the following about this example and made several observations about the data:

- Most user scores are between two and three in terms of their price sensitivity.
- After taking the log value of the user scores, a new user scoring index was created, which recentered the data values around the center of the distribution.
- The pricing scores appear to be stable over time, as the duration of the customer does not seem to have significant influence on the user pricing score. Instead, it appears to be relatively constant over time, within a small band of user scores.

At this point, the analysts may want to explore the range of price tiers offered to customers. Figures 12-22 and 12-23 demonstrate examples of the price tiering currently in place within the customer base.

Figure 12-23 shows the price distribution for a customer base. In this example, loyalty score and price are positively correlated; as the loyalty score increases, so do the prices that the customers are willing to pay. It may seem like a strange phenomenon that the most loyal customers in this example are willing to pay higher prices, but the reality is that customers who are very loyal tend to be less sensitive to price fluctuations or increases. The key, however, is to understand which customers are highly loyal so that appropriate pricing can be charged to the right groups of people.

Figure 12-24 shows a variation on 12-23. In this case, the new graphic portrays the same customer price tiers, but this time a rug representation (Chapter 3) has been added at the bottom to reflect the distribution of the data points.

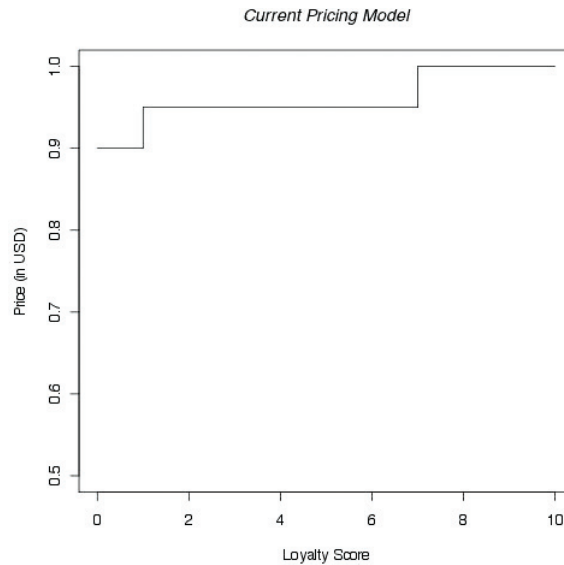


FIGURE 12-23 Graph comparing the price in U.S. dollars with a customer loyalty score

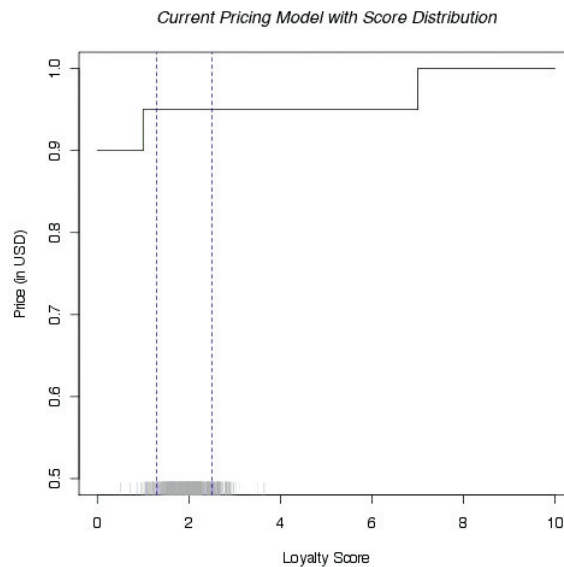


FIGURE 12-24 Graph comparing the price in U.S. dollars with a customer loyalty score (with rug representation)

This rug indicates that the majority of customers in this example are in a tight band of loyalty scores, between about 1 and 3 on the x-axis, all of which offered the same set of prices, which are high (between 0.9 and 1.0 on the y-axis). The y-axis in this example may represent a pricing score, or the raw value of a customer in millions of dollars. The important aspect is to recognize that the pricing is high and is offered consistently to most of the customers in this example.

Based on what was shown in Figure 12-25, the team may decide to develop a new pricing model. Rather than offering static prices to customers regardless of their level of loyalty, a new pricing model might offer more dynamic price points to customers. In this visualization, the data shows the price increases as more of a curvilinear slope relative to the customer loyalty score. The rug at the bottom of the graph indicates that most customers remain between 1 and 3 on the x-axis, but now rather than offering all these customers the same price, the proposal suggests offering progressively higher prices as customer loyalty increases. In one sense, this may seem counterintuitive. It could be argued that the best prices should be offered to the most loyal customers. However, in reality, the opposite is often the case, with the most attractive prices being offered to the least loyal customers. The rationale is that loyal customers are less price sensitive and may enjoy the product and stay with it regardless of small fluctuations in price. Conversely, customers who are not very loyal may defect unless they are offered more attractive prices to stay. In other words, less loyal customers are more price sensitive. To address this issue, a new pricing model that accounts for this may enable an organization to maximize revenue and minimize attrition by offering higher prices to more loyal customers and lower prices to less loyal customers. Creating an iterative depicting the data visually allows the viewer to see these changes in a more concrete way than by looking at tables of numbers or raw values.

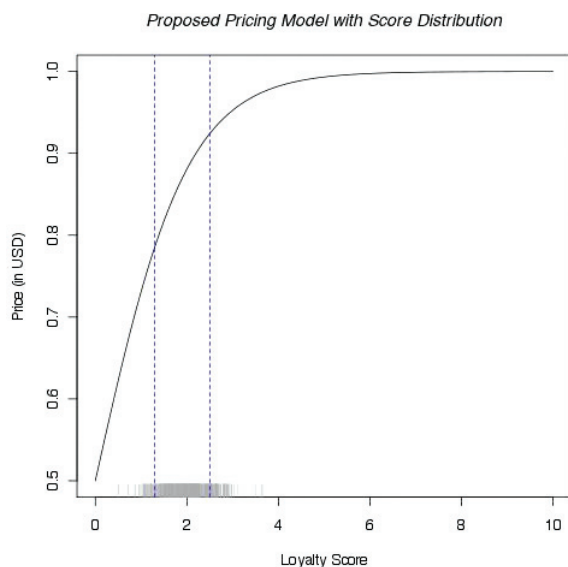


FIGURE 12-25 *New proposed pricing model compared to prices in U.S. dollars with rug*

Data scientists typically iterate and view data in many different ways, framing hypotheses, testing them, and exploring the implications of a given model. This case explores visual examples of pricing distributions, fluctuations in pricing, and the differences in price tiers before and after implementing a new model to optimize price. The visualization work illustrates how the data may look as the result of the model, and helps a data scientist understand the relationships within the data at a glance.

The resulting graph in the pricing scenario appears to be technical regarding the distribution of prices throughout a customer base and would be suitable for a technical audience composed of other data scientists. Figure 12-26 shows an example of how one may present this graphic to an audience of other data scientists or data analysts. This demonstrates a curvilinear relationship between price tiers and customer loyalty when expressed as an index. Note that the comments to the right of the graph relate to the precision of the price targeting, the amount of variability in robustness of the model, and the expectations of model speed when run in a production environment.

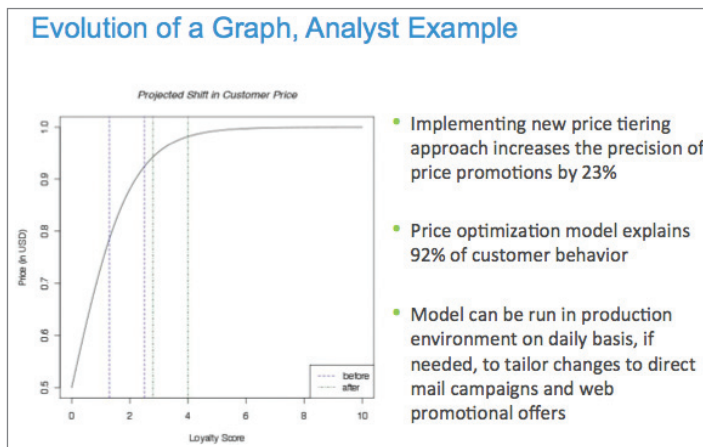


FIGURE 12-26 *Evolution of a graph, analyst example with supporting points*

Figure 12-27 portrays another example of the output from the price optimization project scenario, showing how one may present this to an audience of project sponsors. This demonstrates a simple bar chart depicting the average price per customer or user segment. Figure 12-27 shows a much simpler-looking visual than Figure 12-26. It clearly portrays that customers with lower loyalty scores tend to get lower prices due to targeting from price promotions. Note that the right side of the image focuses on the business impact and cost savings rather than the detailed characteristics of the model.

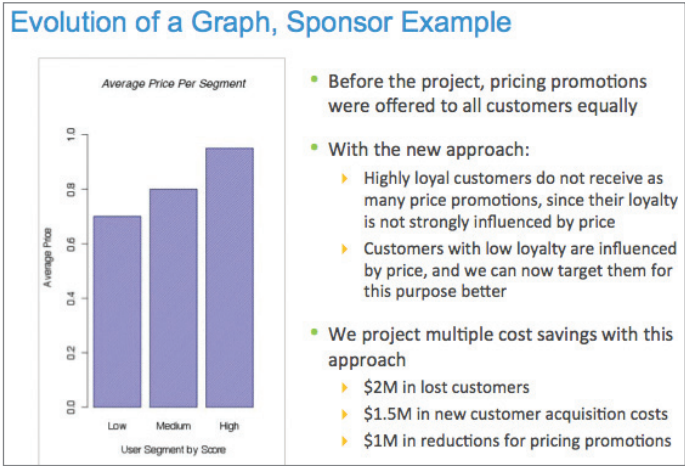


FIGURE 12-27 Evolution of a graph, sponsor example

The comments to the right side of the graphic in Figure 12-27 explain the impact of the model at a high level and the cost savings of implementing this approach to price optimization.

12.3.3 Common Representation Methods

Although there are many types of data visualizations, several fundamental types of charts portray data and information. It is important to know when to use a particular type of chart or graph to express a given kind of data. Table 12-3 shows some basic chart types to guide the reader in understanding that different types of charts are more suited to a situation depending on specific kinds of data and the message the team is attempting to portray. Using a type of chart for data it is not designed for may look interesting or unusual, but it generally confuses the viewer. The objective for the author is to find the best chart for expressing the data clearly so the visual does not impede the message, but rather supports the reader in taking away the intended message.

TABLE 12-3 Common Representation Methods for Data and Charts

Data for Visualization	Type of Chart
Components (parts of whole)	Pie chart
Item	Bar chart
Time series	Line chart
Frequency	Line chart or histogram
Correlation	Scatterplot, side-by-side bar charts

Table 12-3 shows the most fundamental and common data representations, which can be combined, embellished, and made more sophisticated depending on the situation and the audience. It is recommended

that the team consider the message it is trying to communicate and then select the appropriate type of visual to support the point. Misusing charts tends to confuse an audience, so it is important to take into account the data type and desired message when choosing a chart.

Pie charts are designed to show the components, or parts relative to a whole set of things. A pie chart is also the most commonly misused kind of chart. If the situation calls for using a pie chart, employ it only when showing only 2–3 items in a chart, and only for sponsor audiences.

Bar charts and line charts are used much more often and are useful for showing comparisons and trends over time. Even though people use vertical bar charts more often, horizontal bar charts allow an author more room to fit the text labels. Vertical bar charts tend to work well when the labels are small, such as when showing comparisons over time using years.

For frequency, histograms are useful for demonstrating the distribution of data to an analyst audience or to data scientists. As shown in the pricing example earlier in this chapter, data distributions are typically one of the first steps when visualizing data to prepare for model planning. To qualitatively evaluate correlations, scatterplots can be useful to compare relationships among variables.

As with any presentation, consider the audience and level of sophistication when selecting the chart to convey the intended message. These charts are simple examples but can easily become more complex when adding data variables, combining charts, or adding animation where appropriate.

12.3.4 How to Clean Up a Graphic

Many times software packages generate a graphic for a dataset, but the software adds too many things to the graphic. These added visual distractions can make the visual appear busy or otherwise obscure the main points that are to be made with the graphic. In general, it is a best practice to strive for simplicity when creating graphics and data visualization graphs. Knowing how to simplify graphics or clean up a messy chart is helpful for conveying the key message as clearly as possible. Figure 12-28 portrays a line chart with several design problems.

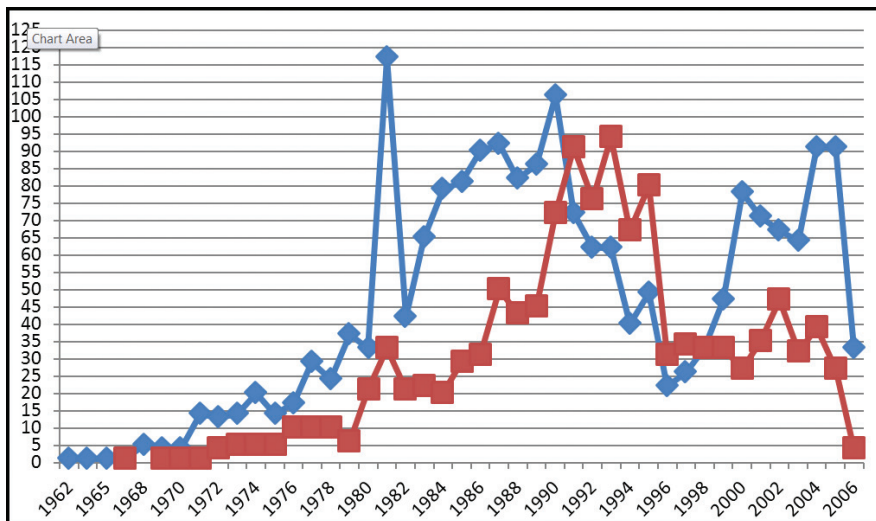


FIGURE 12-28 How to clean up a graphic, example 1 (before)

How to Clean Up a Graphic

The line chart shown in Figure 12-28 compares two trends over time. The chart looks busy and contains a lot of chart junk that distracts the viewer from the main message. *Chart junk* refers to elements of data visualization that provide additional materials but do not contribute to the data portion of the graphic. If chart junk were removed, the meaning and understanding of the graphic would not be diminished; it would instead be made clearer. There are five main kinds of “chart junk” in Figure 12-28:

- **Horizontal grid lines:** These serve no purpose in this graphic. They do not provide additional information for the chart.
- **Chunky data points:** These data points represented as large square blocks draw the viewer’s attention to them but do not represent any specific meaning aside from the data points themselves.
- **Overuse of emphasis colors in the lines and border:** The border of the graphic is a thick, bold line. This forces the viewer’s attention to the perimeter of the graphic, which contains no information value. In addition, the lines showing the trends are relatively thick.
- **No context or labels:** The chart contains no legend to provide context as to what is being shown. The lines also lack labels to explain what they represent.
- **Crowded axis labels:** There are too many axis labels, so they appear crowded. There is no need for labels on the y-axis to appear every five units or for values on the x-axis to appear every two units. Shown in this way, the axis labels distract the viewer from the actual data that is represented by the trend lines in the chart.

The five forms of chart junk in Figure 12-28 are easily corrected, as shown in Figure 12-29. Note that there is no clear message associated with the chart and no legend to provide context for what is shown in Figure 12-28.

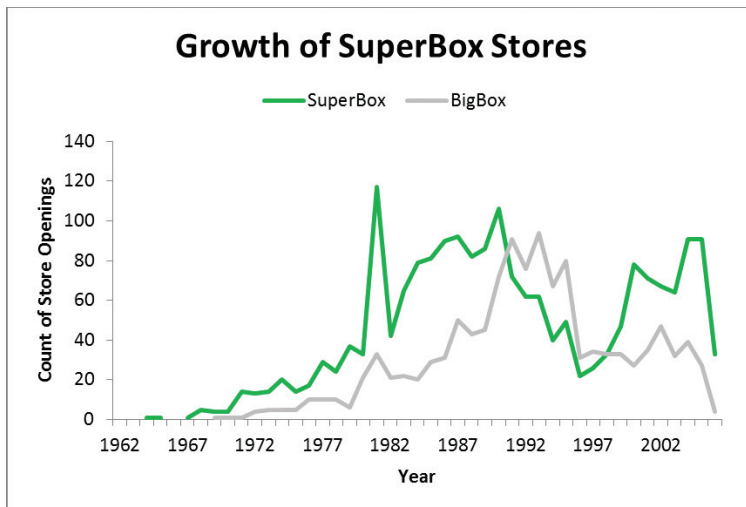


FIGURE 12-29 How to clean up a graphic, example 1 (after)

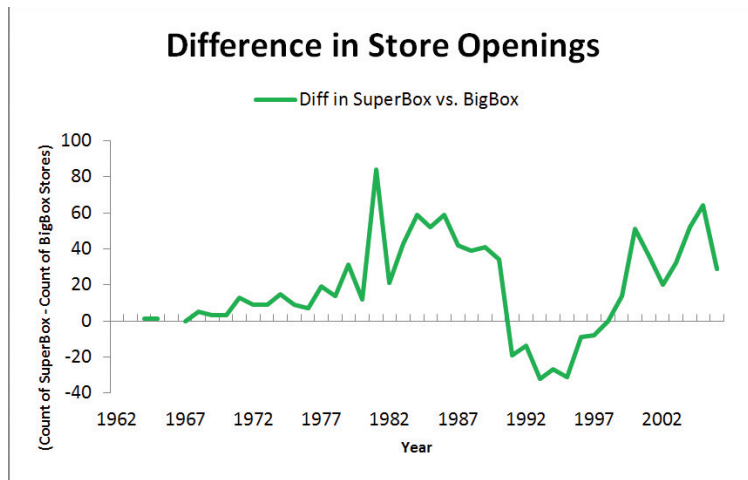


FIGURE 12-30 *How to clean up a graphic, example 1 (alternate “after” view)*

Figures 12-29 and 12-30 portray two examples of cleaned-up versions of the chart shown in Figure 12-28. Note that the problems with chart junk have been addressed. There is a clear label and title for each chart to reinforce the message, and color has been used in ways to highlight the point the author is trying to make. In Figure 12-29, a strong, green color is shown to represent the count of SuperBox stores, because this is where the viewer’s focus should be drawn, whereas the count of BigBox stores is shown in a light gray color.

In addition, note the amount of white space being used in each of the two charts shown in Figures 12-29 and 12-30. Removing grid lines, excessive axes, and the visual noise within the chart allows clear contrast between the emphasis colors (the green line charts) and the standard colors (the lighter gray of the BigBox stores). When creating charts, it is best to draw most of the main visuals in standard colors, light tones, or color shades so that stronger emphasis colors can highlight the main points. In this case, the trend of BigBox stores in light gray fades into the background but does not disappear, while making the SuperBox stores trend in a darker gray (bright green in the online chart) makes it prominent to support the message the author is making about the growth of the SuperBox stores.

An alternative to Figure 12-29 is shown in Figure 12-30. If the main message is to show the difference in the growth of new stores, Figure 12-30 can be created to further simplify Figure 12-28 and graph only the difference between SuperBox stores compared to regular BigBox stores. Two examples are shown to illustrate different ways to convey the message, depending on what it is the author of these charts would like to emphasize.

How to Clean Up a Graphic, Second Example

Another example of cleaning up a chart is portrayed in Figure 12-31. This vertical bar chart suffers from more of the typical problems related to chart junk, including misuse of color schemes and lack of context.

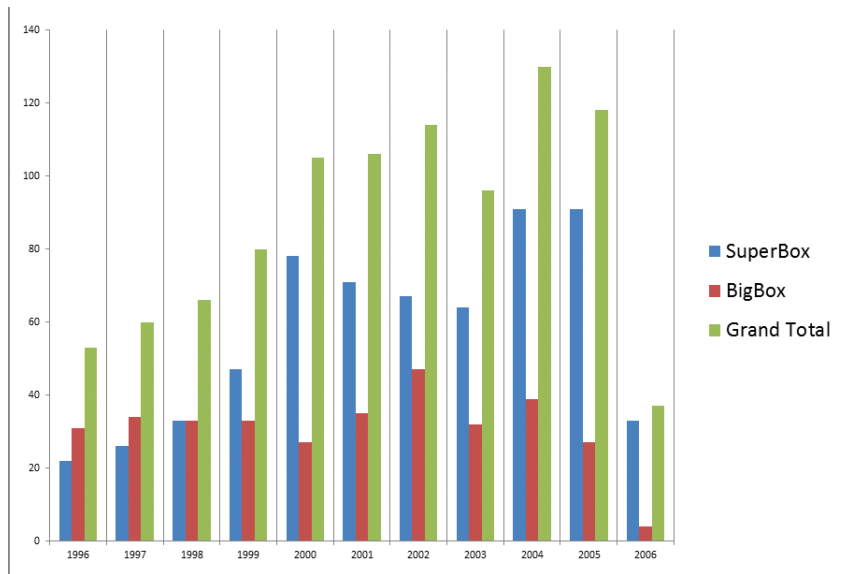


FIGURE 12-31 *How to clean up a graphic, example 2 (before)*

There are five main kinds of chart junk in Figure 12-31:

- **Vertical grid lines:** These vertical grid lines are not needed in this graphic. They provide no additional information to help the viewer understand the message in the data. Instead, these vertical grid lines only distract the viewer from looking at the data.
- **Too much emphasis color:** This bar chart uses strong colors and too much high-contrast dark gray-scale. In general, it is best to use subtle tones, with a low contrast gray as neutral color, and then emphasize the data underscoring the key message in a dark tone or strong color.
- **No chart title:** Because the graphic lacks a chart title, the viewer is not oriented to what he is viewing and does not have proper context.
- **Legend at right restricting chart space:** Although there is a legend for the chart, it is shown on the right side, which causes the vertical bar chart to be compressed horizontally. The legend would make more sense placed across the top, above the chart, where it would not interfere with the data being expressed.
- **Small labels:** The horizontal and vertical axis labels have appropriate spacing, but the font size is too small to be easily read. These should be slightly larger to be easily read, while not appearing too prominent.

Figures 12-32 and 12-33 portray two examples of cleaned-up versions of the chart shown in Figure 12-31. The problems with chart junk have been addressed. There is a clear label and title for each chart to reinforce the message, and appropriate colors have been used in ways to highlight the point the author is trying to make. Figures 12-32 and 12-33 show two options for modifying the graphic, depending on the main point the presenter is trying to make.

Figure 12-32 shows strong emphasis color (dark blue) representing the SuperBox stores to support the chart title: Growth of SuperBox Stores.

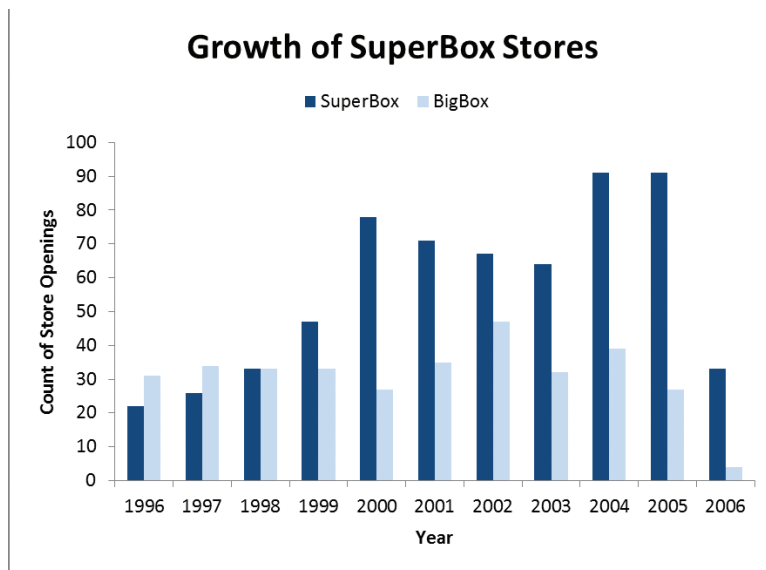


FIGURE 12-32 *How to clean up a graphic, example 2 (after)*

Suppose the presenter wanted to talk about the total growth of BigBox stores instead. A line chart showing the trends over time would be a better choice, as shown in Figure 12-33.

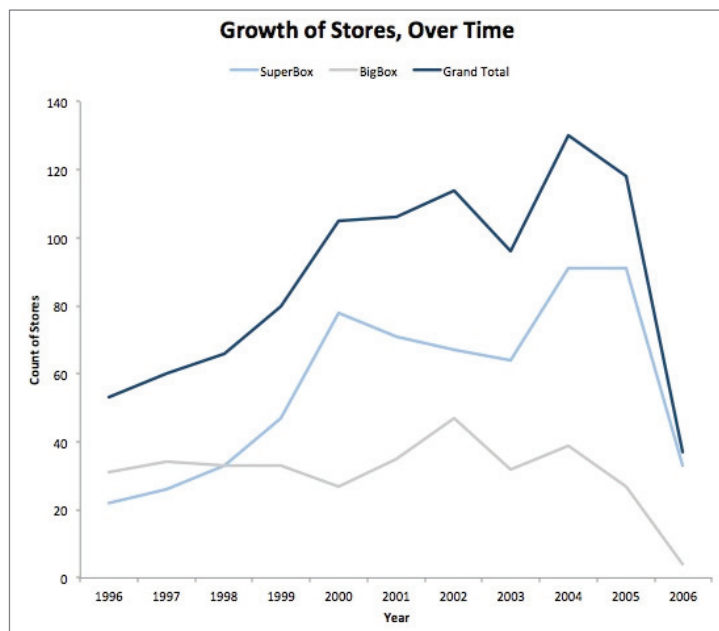


FIGURE 12-33 *How to clean up a graphic, example 2 (alternate view of “after”)*

In both cases, the noise and distractions within the chart have been removed. As a result, the data in the bar chart for providing context has been deemphasized, while other data has been made more prominent because it reinforces the key point as stated in the chart's title.

12.3.5 Additional Considerations

As stated in the previous examples, the emphasis should be on simplicity when creating charts and graphs. Create graphics that are free of chart junk and utilize the simplest method for portraying graphics clearly. The goal of data visualization should be to support the key messages being made as clearly as possible and with few distractions.

Similar to the idea of removing chart junk is being cognizant of the data-ink ratio. *Data-ink* refers to the actual portion of a graphic that portrays the data, while *non-data ink* refers to labels, edges, colors, and other decoration. If one imagined the ink required to print a data visualization on paper, the data-ink ratio could be thought of as (data-ink)/(total ink used to print the graphic). In other words, the greater the ratio of data-ink in the visual, the more data rich it is and the fewer distractions it has [4].

Avoid Using Three-Dimensions in Most Graphics

One more example where people typically err is in adding unnecessary shading, depth, or dimensions to graphics. Figure 12-34 shows a vertical bar chart with two visible dimensions. This example is simple and easy to understand, and the focus is on the data, not the graphics. The author of the chart has chosen to highlight the SuperBox stores in a dark blue color, while the BigBox bars in the chart are in a lighter blue. The title is about the growth of SuperBox stores, and the SuperBox bars in the chart are in a dark, high-contrast shade that draws the viewer's attention to them.

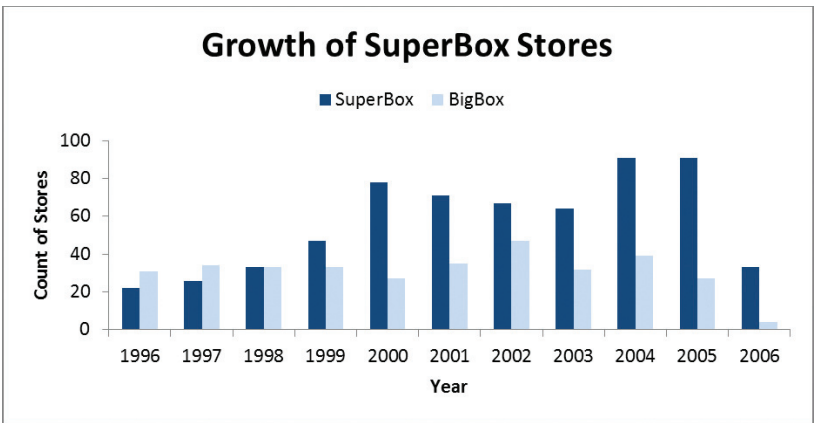


FIGURE 12-34 Simple bar chart, with two dimensions

Compare Figure 12-34 to Figure 12-35, which shows a three-dimensional chart. Figure 12-35 shows the original bar chart at an angle, with some attempt at showing depth. This kind of three-dimensional perspective makes it more difficult for the viewer to gauge the actual data and the scaling becomes deceptive.

Three-dimensional charts often distort scales and axes, and impede viewer cognition. Adding a third dimension for depth in Figure 12-35, does not make it fancier, just more difficult to understand.

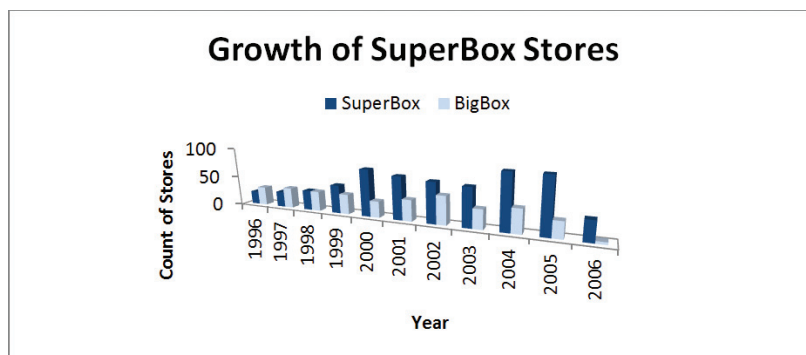


FIGURE 12-35 *Misleading bar chart, with three dimensions*

The charts in Figures 12-34 and 12-35 portray the same data, but it is more difficult to judge the actual height of the bars in Figure 12-35. Moreover, the shadowing and shape of the chart cause most viewers to spend time looking at the perspective of the chart rather than the height of the bars, which is the key message and purpose of this data visualization.

Summary

Communicating the value of analytical projects is critical for sustaining the momentum of a project and building support within organizations. This support is instrumental in turning a successful project into a system or integrating it properly into an existing production environment. Because an analytics project may need to be communicated to audiences with mixed backgrounds, this chapter recommends creating four deliverables to satisfy most of the needs of various stakeholders.

- A presentation for a project sponsor
- A presentation for an analytical audience
- Technical specification documents
- Well-annotated production code

Creating these deliverables enables the analytics project team to communicate and evangelize the work that it did, whereas the code and technical documentation assists the team that wants to implement the models within the production environment.

This chapter illustrates the importance of selecting clear and simple visual representations to support the key points in the final presentations or for portraying data. Most data representations and graphs can be improved by simply removing the visual distractions. This means minimizing or removing chart junk, which distracts the viewer from the main purpose of a chart or graph and does not add information value.

Following several common-sense principles about minimizing distractions in slides and visualizations, communicating clearly and simply, using color in a deliberate way, and taking time to provide context addresses most of the common problems in charts and slides. These few guidelines support the creation of crisp, clear visuals that convey the key messages.

In most cases, the best data visualizations use the simplest, clearest visual to illustrate the key point. Avoid unnecessary embellishment and focus on trying to find the best, simplest method for transmitting the message. Context is critical to orient the viewer to a chart or graph, because people have immediate reactions to imagery on a precognitive level. To this end, make sure to employ thoughtful use of color and orient the viewer with scales, legends, and axes.

Exercises

1. Describe four common deliverables for an analytics project.
2. What is the focus of a presentation for a project sponsor?
3. Give examples of appropriate charts to create in a presentation for other data analysts and data scientists as part of a final presentation. Explain why the charts are appropriate to show each audience.
4. Explain what types of graphs would be appropriate to show data changing over time and why.
5. As part of operationalizing an analytics project, which deliverable would you expect to provide to a Business Intelligence analyst?

References and Further Reading

Following are additional references to learn more about best practices for giving presentations.

- ***Say It with Charts*, by Gene Zelazny [3]:** Simple reference book on how to select the right graphical approach for portraying data and for ensuring the message is clearly conveyed in presentations.
- ***Pyramid Principle*, by Barbara Minto [5]:** Minto pioneered the approach for constructing logical structures for presentations in threes: three sections to the presentations, each with three main points. This teaches people how to weave a story out of the disparate pieces.
- ***Presentation Zen*, by Garr Reynolds [6]:** Teaches how to convey ideas simply and clearly and use imagery in presentations. Shows many before and after versions of graphics and slides.
- ***Now You See It*, by Stephen Few [4]:** Provides many examples for matching the appropriate kind of data visualization to a given dataset.

Bibliography

- [1] N. Yau, "flowingdata.com" [Online]. Available: <http://flowingdata.com>.
- [2] N. Yau, *Visualize This*, Indianapolis: Wiley, 2011.
- [3] G. Zelazny, *Say It with Charts: The Executive's Guide to Visual Communication*, McGraw-Hill, 2001.
- [4] S. Few, *Now You See It: Simple Visualization Techniques for Quantitative Analysis*, Analytics Press, 2009.
- [5] B. Minto, *The Minto Pyramid Principle: Logic in Writing, Thinking, and Problem Solving*, Prentice Hall, 2010.
- [6] G. Reynolds, *Presentation Zen: Simple Ideas on Presentation Design and Delivery*, Berkeley: New Riders, 2011.

Index

Numbers & Symbols

- \ (backward slash) as separator, 69
- / (forward slash) as separator, 69
- 1-itemsets, 147
- 2-itemsets, 148–149
- 3 Vs (volume, variety, velocity), 2–3
- 3-itemsets, 149–150
- 4-itemsets, 150–151

A

- accuracy, 225
- ACF (autocorrelation function), 236–237
- ACME text analysis example, 259–260
 - raw text collection, 260–263
- aggregates (SQL)
 - ordered, 351–352
 - user-defined, 347–351
- aggregators of data, 18
- AIE (Applied Information Economics), 28
- algorithms
 - clustering, 134–135
 - decision trees, 197–200
 - C4.5, 203–204
 - CART, 204
 - ID3, 203
- Alphine Miner, 42
- alternative hypothesis, 102–103
- analytic projects
 - Approach, 369–371
 - BI analyst, 362
 - business users, 361
 - code, 362, 376–377
 - communication, 360–361
 - data engineer, 362
 - data scientists, 362
 - DBA (Database Administrator), 362
 - deliverables, 362–364
 - audiences, 364–365
 - core material, 364–365
 - key points, 372
 - Main Findings, 367–369
 - model description, 371
 - model details, 372–374
 - operationalizing, 360–361
 - outputs, 361
 - presentations, 362
 - Project Goals, 365–367
 - project manager, 362
 - project sponsor, 361
 - recommendations, 374–375
 - stakeholders, 361–362
 - technical specifications, 376–377
- analytic sandboxes. *See* sandboxes
- analytical architecture, 13–15
- analytics
 - business drivers, 11
 - examples, 22–23
 - new approaches, 16–19
- ANOVA, 110–114
- Anscombe's quartet, 82–83
- aov () function, 78
- Apache Hadoop. *See* Hadoop
- APIs (application programming interfaces), Hadoop, 304–305
- apriori () function, 146, 152–157
- Apriori algorithm, 139
 - grocery store example, 143
 - Groceries dataset, 144–146
 - itemset generation, 146–151
 - rule generation, 152–157
 - itemsets, 139, 140–141
 - counting, 158
 - partitioning and, 158
 - sampling and, 158
 - transaction reduction and, 158
- architecture, analytical, 13–15
- arima () function, 246
- ARIMA (Autoregressive Integrated Moving Average) model, 236
 - ACF, 236–237
 - ARMA model, 241–244
 - autoregressive models, 238–239
 - building, 244–252
 - cautions, 252–253
 - constant variance, 250–251
 - evaluating, 244–252
 - fitted time series models, 249–250
 - forecasting, 251–252
 - moving average models, 239–241
 - normality, 250–251
 - PACF, 238–239
 - reasons to choose, 252–253
 - seasonal autoregressive integrated moving average model, 243–244
 - VARIMA, 253
- ARMA (Autoregressive Moving Average) model, 241–244
- array () function, 74
- arrays
 - matrices, 74
 - R, 74–75
- association rules, 138–139
 - application, 143
 - candidate rules, 141–142
 - diagnostics, 158

- testing and, 157–158
- validation, 157–158
- attributes
 - objects, k-means, 130–131
 - R, 71–72
- AUC (area under the curve), 227
- autoregressive models, 238–239
- averages, moving average models, 239–241

B

- bagging, 228
- bag-of-words in text analysis, 265–266
- banking, 18
- `barplot()` function, 88
- barplots, 93–94
- Bayes' Theorem, 212–214. *See also* naïve Bayes
 - conditional probability, 212
- BI (business intelligence)
 - analytical tools, 10
 - versus* Data Science, 12–13
- Big Data
 - 3 Vs, 2–3
 - analytics, examples, 22–23
 - characteristics, 2
 - definitions, 2–3
 - drivers, 15–16
 - ecosystem, 16–19
 - key roles, 19–22
 - McKinsey & Co. on, 3
 - volume, 2–3
- boosting, 228–229
- bootstrap aggregation, 228
- box-and-whisker plots, 95–96
- Box-Jenkins methodology, 235–236
 - ARIMA model, 236
- branches (decision trees), 193
- Brown Corpus, 267–268
- business drivers for analytics, 11
- Business Intelligence Analyst, Operationalize phase, 52
- Business Intelligence Analyst role, 27
- Business User, Operationalize phase, 52
- Business User role, 27
- buyers of data, 18

C

- C4.5 algorithm, 203–204
- cable TV providers, 17
- candidate rules, 141–142
- CART (Classification And Regression Trees), 204

- case folding in text analysis, 264–265
- categorical algorithms, 205
- categorical variables, 170–171
- `cbind()` function, 78
- centroids, 120–122
 - starting positions, 134
- character data types, R, 72
- charts, 386–387
- churn rate (customers), 120
 - logistic regression, 180–181
- `class()` function, 72
- classification
 - bagging, 228
 - boosting, 228–229
 - bootstrap aggregation, 228
 - decision trees, 192–193
 - algorithms, 197–200, 203–204
 - binary decisions, 206
 - branches, 193
 - categorical attributes, 205
 - classification trees, 193
 - correlated variables, 206
 - decision stump, 194
 - evaluating, 204–206
 - greedy algorithm, 204
 - internal nodes, 193
 - irrelevant variables, 205
 - nodes, 193
 - numerical attributes, 205
 - R and, 206–211
 - redundant variables, 206
 - regions, 205
 - regression trees, 193
 - root, 193
 - short trees, 194
 - splits, 193, 194, 197, 200–203
 - structure, 205
 - uses, 194
 - naïve Bayes, 211–212
 - Bayes' theorem, 212–214
 - diagnostics, 217–218
 - naïve Bayes classifier, 214–217
 - R and, 218–224
 - smoothing, 217
- classification trees, 193
- classifiers
 - accuracy, 225
 - diagnostics, 224–228
 - recall, 225
- clickstream, 9
- clustering, 118
 - algorithms, 134–135
 - centroids, 120–122

- starting positions, 134
- diagnostics, 128–129
- k-means, 118–119
 - algorithm, 120–122
 - customer segmentation, 120
 - image processing and, 119
 - medical uses, 119
 - reasons to choose, 130–134
 - rescaling, 133–134
 - units of measure, 132–133
- labels, 127
- number of clusters, 123–127
- code, technical specifications in project, 376–377
- coefficients, linear regression, 169
- combiners, 302–303
- Communicate Results phase of lifecycle, 30, 49–50
- components, short trees as, 194
- conditional entropy, 199
- conditional probability, 212
 - naïve Bayes classifier, 215–216
- confidence, 141–142
 - outcome, 172
 - parameters, 171
- confidence interval, 107
- `confint()` function, 171
- confusion matrix, 224, 280
- contingency tables, 79
- continuous variables, discretization, 211
- corpora
 - Brown Corpus, 267–268
 - corpora in Natural Language Processing, 256
 - IC (information content), 268–269
 - sentiment analysis and, 278
- correlated variables, 206
- credit card companies, 2
- CRISP-DM, 28
- crowdsourcing, 17
- CSV (comma-separated-value) files, 64–65
 - importing, 64–65
- customer segmentation
 - k-means, 120
 - logistic regression, 180–181
- CVS files, 6
- cyclic components of time series analysis, 235

D

- data
 - growth needs, 9–10
 - sources, 15–16
- `data()` function, 84
- data aggregators, 17–18

- data analysis, exploratory, 80–82
 - visualization and, 82–85
- Data Analytics Lifecycle
 - Business Intelligence Analyst role, 27
 - Business User role, 27
 - Communicate Results phase, 30, 49–50
 - GINA case study, 58–59
 - Data Engineer role, 27–28
 - Data preparation phase, 29, 36–37
 - Alpine Miner, 42
 - data conditioning, 40–41
 - data visualization, 41–42
 - Data Wrangler, 42
 - dataset inventory, 39–40
 - ETLT, 38–39
 - GINA case study, 55–56
 - Hadoop, 42
 - OpenRefine, 42
 - sandbox preparation, 37–38
 - tools, 42
 - Data Scientist role, 28
- DBA (Database Administrator) role, 27
- Discovery phase, 29
 - business domain, 30–31
 - data source identification, 35–36
 - framing, 32–33
 - GINA case study, 54–55
 - hypothesis development, 35
 - resources, 31–32
 - sponsor interview, 33–34
 - stakeholder identification, 33
- GINA case study, 53–60
- Model Building phase, 30, 46–48
 - Alpine Miner, 48
 - GINA case study, 56–58
 - Mathematica, 48
 - Matlab, 48
 - Octave, 48
 - PL/R, 48
 - Python, 48
 - R, 48
 - SAS Enterprise Miner, 48
 - SPSS Modeler, 48
 - SQL, 48
 - STATISTICA, 48
 - WEKA, 48
- Model Planning phase, 29–30, 42–44
 - data exploration, 44–45
 - GINA case study, 56
 - model selection, 45
 - R, 45–46

- SAS/ACCESS, 46
 - SQL Analysis services, 46
 - variable selection, 44–45
- Operationalize phase, 30, 50–53, 360
 - Business Intelligence Analyst and, 52
 - Business User and, 52
 - Data Engineer and, 52
 - Data Scientist and, 52
 - DBA (Database Administrator) and, 52
 - GINA case study, 59–60
 - Project Manager and, 52
 - Project Sponsor and, 52
- processes, 28
- Project Manager role, 27
- Project Sponsor role, 27
- roles, 26–28
- data buyers, 18
- data cleansing, 86
- data collectors, 17
- data conditioning, 40–41
- data creation rate, 3
- data devices, 17
- Data Engineer, Operationalize phase, 52
- Data Engineer role, 27–28
- data formats, text analysis, 257
- data frames, 75–76
- data marts, 10
- Data preparation phase of lifecycle, 29, 36–37
 - data conditioning, 40–41
 - data visualization, 41–42
 - dataset inventory, 39–40
 - ETLT, 38–39
 - sandbox preparation, 37–38
- data repositories, 9–11
 - types, 10–11
- Data Savvy Professionals, 20
- Data Science *versus* BI, 12–13
- Data Scientists, 28
 - activities, 20–21
 - business challenges, 20
 - characteristics, 21–22
 - Operationalize phase and, 52
 - recommendations and, 21
 - statistical models and, 20–21
- data sources
 - Discovery phase, 35–36
 - text analysis, 257
- data structures, 5–9
 - quasi-structured data, 6, 7
 - semi-structured data, 6
 - structured data, 6
 - unstructured data, 6
- data types in R, 71–72
 - character, 72
 - logical, 72
 - numeric, 72
 - vectors, 73–74
- data users, 18
- data visualization, 41–42, 377–378
 - CSS and, 378
 - GGobi, 377–378
 - Gnuplot, 377–378
 - graphs, 380–386
 - clean up, 387–392
 - three-dimensional, 392–393
 - HTML and, 378
 - key points with support, 378–379
 - representation methods, 386–387
 - SVG and, 378
- data warehouses, 11
- Data Wrangler, 42
- datasets
 - exporting, R and, 69–71
 - importing, R and, 69–71
 - inventory, 39–40
- Davenport, Tom, 28
- DBA (Database Administrator), 10, 27
 - Operational phase and, 52
- decision trees, 192–193
 - algorithms, 197–200
 - C4.5, 203–204
 - CART, 204
 - categorical, 205
 - greedy, 204
 - ID3, 203
 - numerical, 205
 - binary decisions, 206
 - branches, 193
 - classification trees, 193
 - correlated variables, 206
 - evaluating, 204–206
 - greedy algorithms, 204
 - internal nodes, 193
 - irrelevant variables, 205
 - nodes
 - depth, 193
 - leaf, 193
 - R and, 206–211
 - redundant variables, 206
 - regions, 205
 - regression trees, 193
 - root, 193
 - short trees, 194
 - decision stump, 194

- splits, 193, 197
 - detecting, 200–203
 - limiting, 194
- structure, 205
- uses, 194
- Deep Analytical Talent, 19–20
- DELTA framework, 28
- demand forecasting, linear regression and, 162
- density plots, exploratory data analysis, 88–91
- dependent variables, 162
- descriptive statistics, 79–80
- deviance, 183–184
- devices, 17
 - mobile, 16
 - nontraditional, 16
 - smart devices, 16
- DF (document frequency), 271–272
- diagnostic imaging, 16
- diagnostics
 - association rules, 158
 - classifiers, 224–228
 - linear regression
 - linearity assumption, 173
 - N-fold cross-validation, 177–178
 - normality assumption, 174–177
 - residuals, 173–174
 - logistic regression
 - deviance, 183–184
 - histogram of probabilities, 188
 - log-likelihood test, 184–185
 - pseudo- R^2 , 183
 - ROC curve, 185–187
 - naïve Bayes, 217–218
- `diff()` function, 245
- difference in means, 104
 - confidence interval, 107
 - student's t-testing, 104–106
 - Welch's t-test, 106–108
- differencing, 241–242
- dirty data, 85–87
- Discovery phase of lifecycle, 29
 - data source identification, 35–36
 - framing, 32–33
 - hypothesis development, 35
 - sponsor interview, 33–34
 - stakeholder identification, 33
- discretization of continuous variables, 211
- documents, categorization, 274–277
- `dotchart()` function, 88

E

- Eclipse, 304
- ecosystem of Big Data, 16–19

- Data Savvy Professionals, 20
- Deep Analytical Talent, 19–20
 - key roles, 19–22
 - Technology and Data Enablers, 20
- EDWs (Enterprise Data Warehouses), 10
- effect size, 110
- EMC Google search example, 7–9
- emoticons, 282
- engineering, logistic regression and, 179
- ensemble methods, decision trees, 194
- error distribution
 - linear regression model, 165–166
 - residual standard error, 170
- ETLT, 38–39
- EXCEPT operator (SQL), 333–3334
- exploratory data analysis, 80–82
 - density plot, 88–91
 - dirty data, 85–87
 - histograms, 88–91
 - multiple variables, 91–92
 - analysis over time, 99
 - barplots, 93–94
 - box-and-whisker plots, 95–96
 - dotcharts, 93–94
 - hexbinplots, 96–97
 - versus presentation, 99–101
 - scatterplot matrix, 97–99
 - visualization and, 82–85
 - single variable, 88–91
- exporting datasets in R, 69–71
- expressions, regular, 263

F

- Facebook, 2, 3–4
- factors, 77–78
- financial information, logistic regression and, 179
- FNR (false negative rate), 225
- forecasting
 - ARIMA (Autoregressive Integrated Moving Average)
 - model, 251–252
 - linear regression and, 162
- FP (false positives), confusion matrix, 224
- FPR (false positive rate), 225
- framing in Discovery phase, 32–33
- functions
 - `aov()`, 78
 - `apriori()`, 146, 152–157
 - `arima()`, 246
 - `array()`, 74
 - `barplot()`, 88
 - `cbind()`, 78
 - `class()`, 72
 - `confint()`, 171

- `data()`, 84
- `diff()`, 245
- `dotchart()`, 88
- `gl()`, 84
- `glm()`, 183
- `hclust()`, 135
- `head()`, 65
- `inspect()`, 147, 154–155
- `integer()`, 72
- `IQR()`, 80
- `is.data.frame()`, 75
- `is.na()`, 86
- `is.vector()`, 73
- `jpeg()`, 71
- `kmeans()`, 134
- `kmode()`, 134–135
- `length()`, 72
- `library()`, 70
- `lm()`, 66
- `load.image()`, 68–69
- `matrix.inverse()`, 74
- `mean()`, 86
- `my_range()`, 80
- `na.exclude()`, 86
- `pamk()`, 135
- `Pig`, 307–308
- `plot()`, 65, 153–154, 245
- `predict()`, 172
- `rbind()`, 78
- `read.csv()`, 64–65, 75
- `read.csv2()`, 70
- `read.delim2()`, 70
- `rpart`, 207
- `SQL`, 347–351
- `sqlQuery()`, 70
- `str()`, 75
- `summary()`, 65, 66–67, 79, 80–82
- `t()`, 74
- `ts()`, 245
- `typeof()`, 72
- `wilcox.test()`, 109
- `window functions (SQL)`, 343–347
- `write.csv()`, 70
- `write.csv2()`, 70
- `write.table()`, 70

G

- Generalized Linear Model function, 182
- genetic sequencing, 3, 4
- genomics, 4, 16
- genotyping, 4
- GGobi, 377–378

- GINA (Global Innovation Network and Analysis), Data Analytics Lifecycle case study, 53–60
- `gl()` function, 84
- `glm()` function, 183
- Gnuplot, 377–378
- GPS systems, 16
- Graph Search (Facebook), 3–4
- graphs, 380–386
 - clean up, 387–392
 - three-dimensional, 392–393
- greedy algorithms, 204
- Green Eggs and Ham*, text analysis and, 256
- grocery store example of Apriori algorithm, 143
 - Groceries dataset, 144–146
 - itemsets, frequent generation, 146–151
 - rules, generating, 152–157
- growth needs of data, 9–10
- GUIs (graphical user interfaces), R and, 67–69

H

- Hadoop
 - Data preparation phase, 42
 - Hadoop Streaming API, 304–305
 - HBase, 311–312
 - architecture, 312–317
 - column family names, 319
 - column qualifier names, 319
 - data model, 312–317
 - Java API and, 319
 - rows, 319
 - use cases, 317–319
 - versioning, 319
 - Zookeeper, 319
- HDFS, 300–301
- Hive, 308–311
- LinkedIn, 297
- Mahout, 319–320
- MapReduce, 22
 - combiners, 302–303
 - development, 304–305
 - drivers, 301
 - execution, 304–305
 - mappers, 301–302
 - partitioners, 304
 - structuring, 301–304
- natural language processing, 18
- `Pig`, 306–308
- pipes, 305
- Watson (IBM), 297
- Yahoo!, 297–298
- YARN (Yet Another Resource Negotiator), 305
- hash-based itemsets, Apriori algorithm and, 158

HAWQ (Hadoop With Query), 321
 HBase, 311–312

- architecture, 312–317
- column family names, 319
- column qualifier names, 319
- data model, 312–317
- Java API and, 319
- rows, 319
- use cases, 317–319
- versioning, 319
- Zookeeper, 319

 hclust () function, 135
 HDFS (Hadoop Distributed File System), 300–301
 head () function, 65
 hexbinplots, 96–97
 histograms

- exploratory data analysis, 88–91
- logistic regression, 188

 Hive, 308–311
 HiveQL (Hive Query Language), 308
 Hopper, Grace, 299
 Hubbard, Doug, 28
 HVE (Hadoop Virtualization Extensions), 321
 hypotheses

- alternative hypothesis, 102–103
- Discovery phase, 35
- null hypothesis, 102

 hypothesis testing, 102–104

- two-sided hypothesis testing, 105
- type I errors, 109–110
- type II errors, 109–110

I

IBM Watson, 297
 ID3 algorithm, 203
 IDE (Interactive Development Environment), 304
 IDF (inverted document frequency), 271–272
 importing datasets in R, 69–71
 in-database analytics

- SQL, 328–338
- text analysis, 338–339

 independent variables, 162
 input variables, 192
 inspect () function, 147, 154–155
 integer () function, 72
 internal nodes (decision trees), 193
 Internet of Things, 17–18
 INTERSECT operator (SQL), 333
 IQR () function, 80
 is.data.frame () function, 75
 is.na () function, 86

is.vector () function, 73
 itemsets, 139

- 1-itemsets, 147
- 2-itemsets, 148–149
- 3-itemsets, 149–150
- 4-itemsets, 150–151
- Apriori algorithm, 139
- Apriori property, 139
- downward closure property, 139
- dynamic counting, Apriori algorithm and, 158
- frequent itemset, 139
- generation, frequent, 146–151
- hash-based, Apriori algorithm and, 158
- k-itemset, 139, 140–141

J

joins (SQL), 330–332
 jpeg () function, 71

K

k clusters

- finding, 120–122
- number of, 123–127

 k-itemset, 139, 140–141
 k-means, 118–119

- customer segmentation, 120
- image processing and, 119
- k clusters
 - finding, 120–122
 - number of, 123–127
- medical uses, 119
- objects, attributes, 130–131
- R and, 123–127
- reasons to choose, 130–134
- rescaling, 133–134
- units of measure, 132–133

 kmeans () function, 134
 kmode () function, 134–135

L

lag, 237
 Laplace smoothing, 217
 lasso regression, 189
 LDA (latent Dirichlet allocation), 274–275
 leaf nodes, 192, 193
 lemmatization, text analysis and, 258
 length () function, 72
 leverage, 142
 library () function, 70

lifecycle. *See also* Data Analytics Lifecycle
lift, 142

linear regression, 162
 coefficients, 169
 diagnostics
 linearity assumption, 173
 N-fold cross-validation, 177–178
 normality assumption, 174–177
 residuals, 173–174
 model, 163–165
 categorical variables, 170–171
 normally distributed errors, 165–166
 outcome confidence intervals, 172
 parameter confidence intervals, 171
 prediction interval on outcome, 172
 R, 166–170
 p-values, 169–170
 use cases, 162–163

LinkedIn, 2, 22–23, 297

lists in R, 76–77

`lm()` function, 66

`load.image()` function, 68–69

logical data types, R, 72

logistic regression, 178
 cautions, 188–189
 diagnostics, 181–182
 deviance, 183–184
 histogram of probabilities, 188
 log-likelihood test, 184–185
 pseudo- R^2 , 183
 ROC curve, 185–187
 Generalized Linear Model function, 182
 model, 179–181
 multinomial, 190
 reasons to choose, 188–189
 use cases, 179

log-likelihood test, 184–185

loyalty cards, 17

M

MAD (Magnetic/Agile/Deep) skills, 28, 352–356

MADlib, 352–356

Mahout, 319–320

MapReduce, 22, 298–299
 combiners, 302–303
 development, 304–305
 drivers, 301–302
 execution, 304–305
 mappers, 301–302
 partitioners, 304
 structuring, 301–304
market basket analysis, 139
 association rules, 143

marketing, logistic regression and, 179

master nodes, 301

matrices
 confusion matrix, 224
 R, 74–75
 scatterplot matrices, 97–99

`matrix.inverse()` function, 74

MaxEnt (maximum entropy), 278

McKinsey & Co. definition of Big Data, 3

`mean()` function, 86

medical information, 16
 k-means and, 119
 linear regression and, 162
 logistic regression and, 179

minimum confidence, 141

missing data, 86

mobile devices, 16

mobile phone companies, 2

Model Building phase of lifecycle, 30, 46–48

 Alpine Miner, 48
 Mathematica, 48
 Matlab, 48
 Octave, 48
 PL/R, 48
 Python, 48
 R, 48
 SAS Enterprise Miner, 48
 SPSS Modeler, 48
 SQL, 48
 STATISTICA, 48
 WEKA, 48

Model Planning phase of lifecycle, 29–30, 42–44

 data exploration, 44–45
 model selection, 45
 R, 45–46
 SAS/ACCESS, 46
 SQL Analysis services, 46
 variables, selecting, 44–45

morphological features in text analysis, 266–267

moving average models, 239–241

MPP (massively parallel processing), 5

MTurk (Mechanical Turk), 282

multinomial logistic regression, 190

multivariate time series analysis, 253

`my_range()` function, 80

N

`na.exclude()` function, 86

naïve Bayes, 211–212
 Bayes' theorem, 212–214
 diagnostics, 217–218

- naïve Bayes classifier, 214–217
- R and, 218–224
- sentiment analysis and, 278
- smoothing, 217
- natural language processing, 18
- N-fold cross-validation, 177–178
- NLP (Natural Language Processing), 256
- nodes
 - master, 301
 - worker, 301
- nodes (decision trees), 192
 - depth, 193
 - leaf, 193
 - leaf nodes, 192, 193
- nonparametric tests, 108–109
- nontraditional devices, 16
- normality
 - ARIMA model, 250–251
 - linear regression, 174–177
- normalization, data conditioning, 40–41
- NoSQL, 322–323
- null deviance, 183
- null hypothesis, 102
- numeric data types, R, 72
- numerical algorithms, 205
- numerical underflow, 216–217

O

- objects, k-means, attributes, 130–131
- OLAP (online analytical processing), 6
 - cubes, 10
- OpenRefine, 42
- Operationalize phase of lifecycle, 30, 50–53, 360
 - Business Intelligence Analyst and, 52
 - Business User and, 52
 - Data Engineer and, 52
 - Data Scientist and, 52
 - DBA (Database Administrator) and, 52
 - Project Manager and, 52
 - Project Sponsor and, 52
- operators, subsetting, 75
- outcome
 - confidence intervals, 172
 - prediction interval, 172

P

- PACF (partial autocorrelation function), 238–239
- `pamk ()` function, 135
- parameters, confidence intervals, 171
- parametric tests, 108–109

- parsing, text analysis and, 257
- partitioning
 - Apriori algorithm and, 158
 - MapReduce, 304
- photographs, 16
- Pig, 306–308
- Pivotal HD Enterprise, 320–321
- `plot ()` function, 65, 153–154, 245
- POS (part-of-speech) tagging, 258
- power of a test, 110
- precision in sentiment analysis, 281
- `predict ()` function, 172
- prediction trees. *See* decision trees
- presentation *versus* data exploration, 99–101
- probability, conditional, 212
 - naïve Bayes classifier, 215–216
- Project Manager, Operationalize phase, 52
- Project Manager role, 27
- Project Sponsor, Operationalize phase, 52
- Project Sponsor role, 27
- pseudo- R^2 , 183
- p-values, linear regression, 169–170

Q

- quasi-structured data, 6, 7
- queries, SQL, 329–330
 - nested, 3334
 - subqueries, 3334

R

- arrays, 74–75
- attributes, types, 71–72
- data frames, 75–76
- data types, 71–72
 - character, 72
 - logical, 72
 - numeric, 72
 - vectors, 73–74
- decision trees, 206–211
- descriptive statistics, 79–80
- exploratory data analysis, 80–82
 - density plot, 88–91
 - dirty data, 85–87
 - histograms, 88–91
 - multiple variables, 91–99
 - versus* presentation, 99–101
 - visualization and, 82–85, 88–91
- factors, 77–78
- functions

- aov(), 78
- array(), 74
- barplot(), 88
- cbind(), 78
- class(), 72
- data(), 84
- dotchart(), 88
- gl(), 84
- head(), 65
- import function defaults, 70
- integer(), 72
- IQR(), 80
- is.data.frame(), 75
- is.na(), 86
- is.vector(), 73
- jpeg(), 71
- length(), 72
- library(), 70
- lm(), 66
- load.image(), 68–69
- my_range(), 80
- plot() function, 65
- rbind(), 78
- read.csv(), 65, 75
- read.csv2(), 70
- read.delim(), 69
- read.delim2(), 70
- read.table(), 69
- str(), 75
- summary(), 65, 66–67, 79
- t(), 74
- typeof(), 72
- visualizing single variable, 88
- write.csv(), 70
- write.csv2(), 70
- write.table(), 70
- GUIs, 67–69
- import/export, 69–71
- k-means analysis, 123–127
- linear regression model, 166–170
- lists, 76–77
- matrices, 74–75
- model planning and, 45–46
- naïve Bayes and, 218–224
- operators, subsetting, 75
- overview, 64–67
- statistical techniques, 101–102
 - ANOVA, 110–114
 - difference in means, 104–108
 - effect size, 110
 - hypothesis testing, 102–104
 - power of test, 110
 - sample size, 110
 - type I errors, 109–110
 - type II errors, 109–110
- tables, contingency tables, 79
- R commander GUI, 67
- random components of time series analysis, 235
- Rattle GUI, 67
- raw text
 - collection, 260–263
 - tokenization, 264
- rbind() function, 78
- RDBMS, 6
- read.csv() function, 64–65, 75
- read.csv2() function, 70
- read.delim() function, 69
- read.delim2() function, 70
- read.table() function, 69
- real estate, linear regression and, 162
- recall in sentiment analysis, 281
- redundant variables, 206
- regression
 - lasso, 189
 - linear, 162
 - coefficients, 169
 - diagnostics, 173–178
 - model, 163–172
 - p-values, 169–170
 - use cases, 162–163
 - logistic, 178
 - cautions, 188–189
 - diagnostics, 181–188
 - model, 179–181
 - multinomial logistic, 190
 - reasons to choose, 188–189
 - use cases, 179
 - multinomial logistic, 190
 - ridge, 189
 - variables
 - dependent, 162
 - independent, 162
- regression trees, 193
- regular expressions, 263, 339–340
- relationships, 141
- repositories, 9–11
 - types, 10–11
- representation methods, 386–387
- rescaling, k-means, 133–134
- residual deviance, 183
- residual standard error, 170

- residuals, linear regression, 173–174
- resources, Discovery phase of lifecycle, 31–32
- RFID readers, 16
- ridge regression, 189
- ROC (receiver operating characteristic) curve, 185–187, 225
- roots (decision trees), 193
- `rpart` function, 207
- RStudio GUI, 67–68
- rules
 - association rules, 138–139
 - application, 143
 - candidate rules, 141–142
 - diagnostics, 158
 - testing and, 157–158
 - validation, 157–158
 - generating, grocery store example (Apriori), 152–157

S

- sales, time series analysis and, 234
- sample size, 110
- sampling, Apriori algorithm and, 158
- sandboxes, 10, 11. *See also* work spaces
 - Data preparation phase, 37–38
- SAS/ACCESS, model planning, 46
- scatterplot matrix, 97–99
- scatterplots, 81
 - Anscombe's quartet, 83
 - multiple variables, 91–92
- scientific method, 28
- searches, text analysis and, 257
- seasonal autoregressive integrated moving average model, 243–244
- seasonality components of time series analysis, 235
- seismic processing, 16
- semi-structured data, 6
- SensorNet, 17–18
- sentiment analysis in text analysis, 277–283
 - confusion matrix, 280
 - precision, 281
 - recall, 281
- shopping
 - loyalty cards, 17
 - RFID chips in carts, 17
- short trees, 194
- smart devices, 16
- smartphones, 17
- smoothing, 217
- social media, 3–4
- sources of data, 15–16
- split parts planning, time series analysis and, 234–235
- splits (decision trees), 193
 - detecting, 200–203

- sponsor interview, Discovery phase, 33
- spreadmarts, 10
- spreadsheets, 6, 9, 10
- SQL (Structured Query Language), 328–329
 - aggregates
 - ordered, 351–352
 - user-defined, 347–351
 - EXCEPT operator, 333–3334
 - functions, user-defined, 347–351
 - grouping, 334–338
 - INTERSECT operator, 333
 - joins, 330–332
 - MADlib, 352–356
 - queries, 329–330
 - nested, 3334
 - subqueries, 3334
 - set operations, 332–334
 - UNION ALL operator, 332–333
 - window functions, 343–347
- SQL Analysis services, model planning and, 46
- `sqlQuery ()` function, 70
- stakeholders, Discovery phase of lifecycle, 33
- stationary time series, 236
- statistical techniques, 101–102
 - ANOVA, 110–114
 - difference in means, 104
 - student's t-test, 104–106
 - Welch's t-test, 106–108
 - effect size, 110
 - hypothesis testing, 102–104
 - power of test, 110
 - sample size, 110
 - type I errors, 109–110
 - type II errors, 109–110
 - Wilcoxon rank-sum test, 108–109
- statistics
 - Anscombe's quartet, 82–83
 - descriptive, 79–80
- stemming, text analysis and, 258
- stock trading, time series analysis and, 235
- stop words, 270–271
- `str ()` function, 75
- structured data, 6
- subsetting operators, 75
- `summary ()` function, 65, 66–67, 79, 80–82
- SVM (support vector machines), 278

T

- `t ()` function, 74
- tables, contingency tables, 79
- Target stores, 22
- t-distribution

- ANOVA, 110–114
 - student's t-test, 104–106
 - Welch's t-test, 106–108
 - technical specifications in project, 376–377
 - Technology and Data Enablers, 20
 - testing, association rules and, 157–158
 - text analysis, 256
 - ACME example, 259–263
 - bag-of-words, 265–266
 - corpora, 264–265
 - Brown Corpus, 267–268
 - corpora in Natural Language Processing, 256
 - IC (information corpora), 268–269
 - data formats, 257
 - data sources, 257
 - document categorization, 274–277
 - Green Eggs and Ham*, 256
 - in-database, 338–339
 - lemmatization, 258
 - morphological features, 266–267
 - NLP (Natural Language Processing), 256
 - parsing, 257
 - POS (part-of-speech) tagging, 258
 - raw text, collection, 260–263
 - search and retrieval, 257
 - sentiment analysis, 277–283
 - stemming, 258
 - stop words, 270–271
 - text mining, 257–258
 - TF (term frequency) of words, 265–266
 - DF, 271–272
 - IDF, 271–272
 - lemmatization, 271
 - stemming, 271
 - stop words, 270–271
 - TFIDF, 269–274
 - tokenization, 264
 - topic modeling, 267, 274
 - LDA (latent Dirichlet allocation), 274–275
 - web scraper, 262–263
 - word clouds, 284
 - Zipf's Law, 265–266
 - text mining, 257
 - textual data files, 6
 - TF (term frequency) of words, 265–266
 - DF (document frequency), 271–272
 - IDF (inverted document frequency), 271–272
 - lemmatization, 271
 - stemming, 271
 - stop words, 270–271
 - TFIDF, 269–274
 - TFIDF (Term Frequency-Inverse Document Frequency), 269–274, 285–286
 - time series analysis
 - ARIMA model, 236
 - ACF, 236–237
 - ARMA model, 241–244
 - autoregressive models, 238–239
 - building, 244–252
 - cautions, 252–253
 - constant variance, 250–251
 - evaluating, 244–252
 - fitted models, 249–250
 - forecasting, 251–252
 - moving average models, 239–241
 - normality, 250–251
 - PACF, 238–239
 - reasons to choose, 252–253
 - seasonal autoregressive integrated moving average model, 243–244
 - ARMAX (Autoregressive Moving Average with Exogenous inputs), 253
 - Box-Jenkins methodology, 235–236
 - cyclic components, 235
 - differencing, 241–242
 - fitted models, 249–250
 - GARCH (Generalized Autoregressive Conditionally Heteroscedastic), 253
 - Kalman filtering, 253
 - multivariate time series analysis, 253
 - random components, 235
 - seasonal autoregressive integrated moving average model, 243–244
 - seasonality, 235
 - spectral analysis, 253
 - stationary time series, 236
 - trends, 235
 - use cases, 234–235
 - white noise process, 239
 - tokenization in text analysis, 264
 - topic modeling in text analysis, 267, 274
 - LDA (latent Dirichlet allocation), 274–275
 - TP (true positives), confusion matrix, 224
 - TPR (true positive rate), 225
 - transaction data, 6
 - transaction reduction, Apriori algorithm and, 158
 - trends, time series analysis, 235
 - TRP (True Positive Rate), 185–187
 - `ts()` function, 245
 - two-sided hypothesis test, 105
 - type I errors, 109–110
 - type II errors, 109–110
 - `typeof()` function, 72
- ## U
- UNION ALL operator (SQL), 332–333
 - units of measure, k-means, 132–133
 - unstructured data, 6

- Apache Hadoop, HDFS, 300–301
- LinkedIn, 297
- MapReduce, 298–299
- natural language processing, 18
- use cases, 296–298
- Watson (IBM), 297
- Yahoo!, 297–298

unsupervised techniques. *See* clustering
users of data, 18

V

validation, association rules and, 157–158
variables

- categorical, 170–171
- continuous, discretization, 211
- correlated, 206
- decision trees, 205
- dependent, 162
- factors, 77–78
- independent, 162
- input, 192
- redundant, 206

VARIMA (Vector ARIMA), 253

vectors, R, 73–74

video footage, 16

- k-means and, 119

video surveillance, 16

visualization, 41–42. *See also* data visualization

- exploratory data analysis, 82–85

- single variable, 88–91

grocery store example (Apriori), 152–157

volume, variety, velocity. *See* 3 Vs (volume, variety, velocity)

W

Watson (IBM), 297

web scraper, 262–263

white noise process, 239

Wilcoxon rank-sum test, 108–109

`wilcox.test()` function, 109

window functions (SQL), 343–347

word clouds, 284

work spaces, 10, 11. *See also* sandboxes

- Data preparation phase, 37–38

worker nodes, 301

`write.csv()` function, 70

`write.csv2()` function, 70

`write.table()` function, 70

WSS (Within Sum of Squares), 123–127

X-Z

XML (eXtensible Markup Language), 6

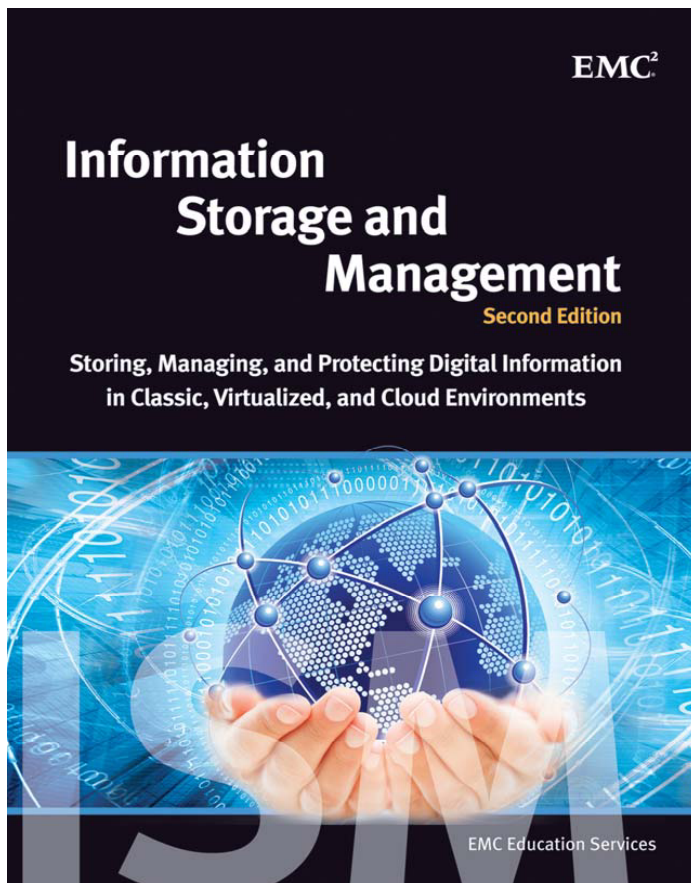
Yahoo!, 297–298

YARN (Yet Another Resource Negotiator), 305

Zipf's Law, 265–266

A comprehensive book on information storage and management by the EMC's Education Services, a global technology leader.

More than ever, the IT industry is challenged with employing and developing highly skilled technical professionals with storage technology expertise across classic, virtualized, and cloud environments. This book covers concepts, principles, and deployment considerations across technologies that are used for storing and managing information.



Key Technology Strategies for Classic, Virtualized, and Cloud Environments:

- Challenges and Solutions for Data Storage and Management
- Intelligent Storage, Object-Based Storage, and Unified Storage
- Storage Networking, Federation, and Protocols
- Backup, Recovery, Deduplication, and Archive
- Business Continuity and Disaster Recovery
- Cloud Computing and Converged Infrastructure
- Storage Security and Managing Storage Infrastructure

978-1-118-09483-9 · \$70.00 US · \$77.00 CAN · £47.50 UK

WILEY

WILEY END USER LICENSE AGREEMENT

Go to www.wiley.com/go/eula to access Wiley's ebook EULA.