

## Chapter 5

# Data Visualization

*The drawing shows me at one glance what might be spread over ten pages in a book.*

— Ivan Turgenev, *Fathers and Sons*

**Abstract** A visual is successful when the information encoded in the data is efficiently transmitted to an audience. Data visualization is the discipline dedicated to the principles and methods of translating data to visual form. In this chapter we discuss the principles that produce successful visualizations. The second section illustrates the principles through examples of best and worst practices. In the final section, we navigate through the construction of our best-example graphics.

### 5.1 Introduction

Humans are visual animals. We absorb sensory information most efficiently through vision. It's no surprise that data visualization is very effective for extracting information from data. As we see throughout this volume, gathering large amounts of data has never been easier. Even 40 years ago, displaying that information was difficult to do well, requiring specialized tools or a steady hand. The democratization of creating figures from data allows us to create more visualizations than ever before. And with this profusion comes the ability to develop best practices. Figures *encode* information from a data set, displaying those data as ink on a page or, more commonly, as pixels on a screen. This encoding makes use of a vernacular that has developed over the last several centuries. The typical audience will understand the Cartesian plane and its axes. We understand how to determine values of points in a scatterplot and easily manage the color-coding of groups.

A visualization is effective when it can be quickly and accurately *decoded* by the audience—the salient points should be almost immediately apparent. Edward Tufte, a vocal evangelist for better graphics, calls this the “Interoocular Trauma Test”: does the visualization hit one immediately between the eyes? A common expression about scripting languages, like `Python` and `Perl`, is that they should make easy things easy and difficult things possible. The same can be said of good visualizations: the key features of the narrative should emerge immediately and the more subtle relationships should be visible whenever possible.

Data visualization combines several different threads and we will cover each in a section. The first is an understanding of the guiding principles of graphics. In this section we lean heavily on the pioneering work of William Cleveland, the researcher who truly brought graphics into the modern era. Most of the excellent graphics one sees in data science journalism<sup>1</sup> are built using the ideas he introduced. The second is a basic understanding of the paradigms of graphics that are often employed for data of different types. Knowing these will allow us to create an abstract version of the graphic. This ability, to sketch a version of the graphic you wish to create, is critical. It is impossible to follow a map if you don’t know where you’re going. Finally, we must be able to tell a software package how to render the image that is in our minds. The best tool for creating data visualizations in the context of an analysis<sup>2</sup> is Hadley Wickham’s `ggplot2` available in R [65, 66]. The “gg” in the package name stands for “The Grammar of Graphics”. You can think of this grammar as being a semi-secret language that, for better or worse, you must be fluent in to realize the potential of graphics from data. Recognizing this reality, we will introduce the grammar and illustrate its implementation in R syntax.

The data used in this chapter and Chap. 10 originates from the largest cooperative (co-op) grocery store in the United States. We received approximately 20 gigabytes of transaction-level data covering 6 years of store activity. The data are essentially cash register or point-of-sales receipts. As is common with automatically recorded point-of-sales information, a considerable amount of associated meta-data is attached to the receipt. We work a great deal with two variables: the department classification of the item (e.g., produce) and whether the shopper is a member of the co-op. The co-op is member-owned and approximately three-quarters of the transaction records

---

<sup>1</sup> Three great examples:

- The Upshot from the New York Times: <http://www.nytimes.com/section/upshot>.
- Five Thirty Eight, Nate Silver’s organization that has largely invented the field of data science journalism. <http://fivethirtyeight.com>.
- Flowing Data, a site created by Nathan Yau dedicated to creating beautiful and informative data visualizations. <http://flowingdata.com>.

<sup>2</sup> If you are building interactive graphics or large-scale graphics via the web, there are better tools. Check out `bootstrap`, `D3`, and `crossfilter`.

originated from members. The remainder originated from non-members. If the shopper is a member, then a unique, anonymous identifier is attached to receipt that allows us to analyze data at the shopper level.

We use the co-op data in this chapter to illustrate a variety of data visualizations. In Chap. 10, we develop a prediction function that classifies non-member shoppers to customer segment with the ultimate goal of better understanding the non-member shoppers.

## 5.2 Principles of Data Visualization

When describing what makes good data visualization, there are two paths to follow: that of simplicity and that of exhaustiveness. There are entire books of the latter variety and thus we opt for a treatment that will give the reader minimal guidelines and that points them in the direction of the more thorough treatments. Our goal, after all, is to be able to make good visualizations and improve on those found in the wild.

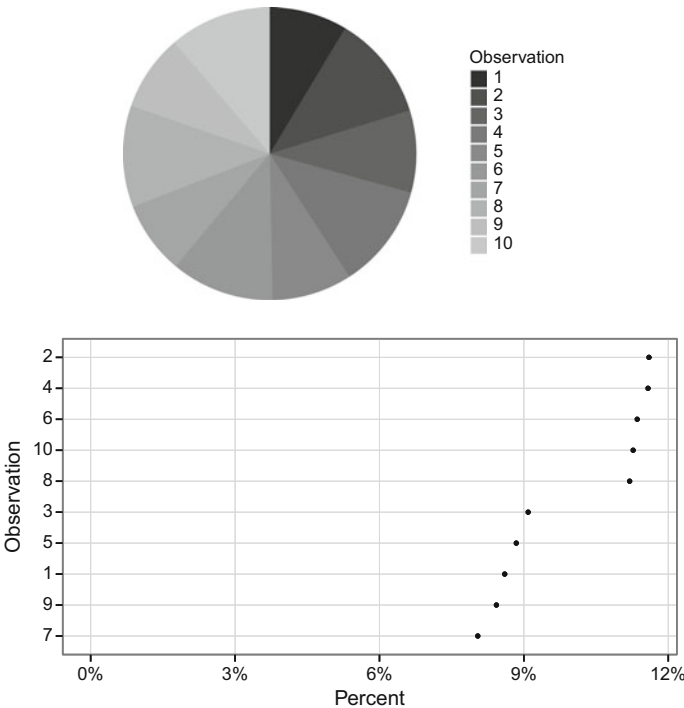
A data visualization is effective when the information that the analyst is trying to convey to the audience is transmitted. Complicated information sometimes requires complicated visualizations, but often we can organize our thoughts in terms of principles.

**Let the Data Speak.** If there is one overarching goal, then this is it: the data should be allowed to speak for itself. As Strunk and White say, “vigorous writing is concise.” Similarly, good data visualization allows the data to tell its story. Tufte coined a term for elements of a visualization that do not add to understanding: *chartjunk*. A good, revised definition of chartjunk from Robert Kosara at Tableau is this: any element of a chart that does not contribute to clarifying the intended message. Historically, the most egregious violations of this principle came from Excel. The default settings in Excel now avoid the worst examples of chartjunk, but options to add them abound, particularly with the addition of mysterious third dimensions to one- or two-dimensional data sets, color coding for no reason, and patterns that impede understanding. For each element of a figure, we must ask if that element is serving our principal goal of letting the data speak.

**Let the Data Speak Clearly.** A subtle addition to our previous point is to let the data tell their story clearly and quickly. As we shall see, the relationship between two variables, often plotted against one another in a scatterplot, can be illuminated by the addition of a smoothed or fitted line. At the other end of the spectrum, Fig. 5.1 illustrates how certain visualizations, in this case the much- and rightly-maligned pie chart can obfuscate the story. In the pie chart version, it’s extremely difficult to identify the predominant pattern—the presence of two sets of observations with different means. The lower panel in Fig. 5.1 shows a dotchart which

lets the data speak for itself and uses a sensible ordering of observations. This visualization conveys much more information and allows for it to be immediately decoded.

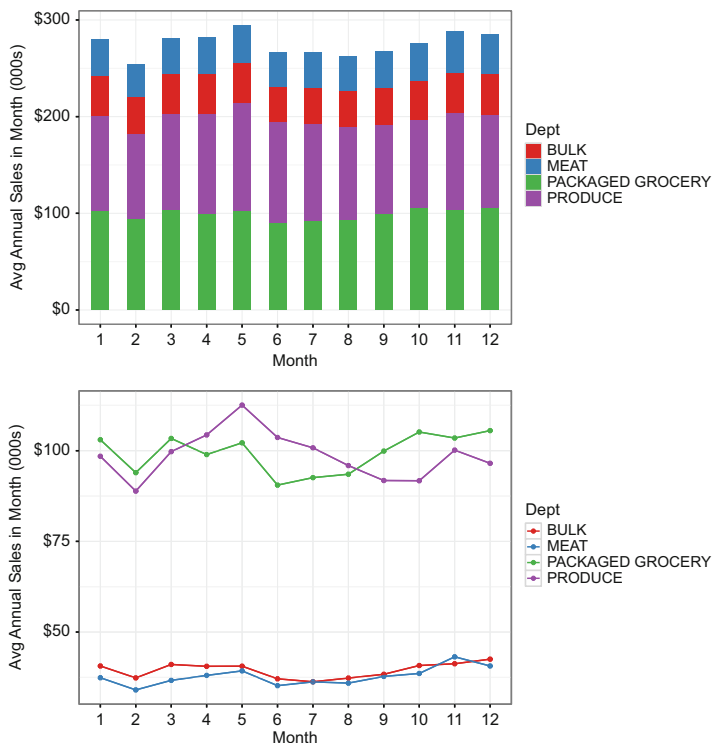
**Fig. 5.1** A pie chart makes patterns in the data difficult to decode; the dotchart is an improvement



We introduce the dotchart in Sect. 5.3 and it is an excellent choice for displaying univariate data. There are other types of graphs besides the pie chart that inhibit understanding, notably stacked bar charts and area charts. The stacked bar chart makes it difficult to understand the behavior of the individual elements, as we see in Fig. 5.2. This figure has two views of average sales data by month for a grocery store for four large departments. In the upper panel, we can see that May is the month of maximum overall sales (by looking at the heights of bars and ignoring the colors). We can see that packaged grocery appears to have lower sales around July and that meat and bulk appear to be smaller departments, although it is difficult to gauge the magnitude of the difference. In the lower panel we have replaced the stacked bars with points connected by lines. The points allow precise estimation of individual observations and the lines (and the color) help us group the departments together across time. Now we see that packaged grocery and produce are larger than the other two departments, by a factor of just more than two. Note that the lower panel figure is less than perfect

because the vertical axis does not extend to 0. We also see that produce has an annual cycle out of phase with the other departments, peaking in the North American summer months. Area charts such as the pie chart typically arise when circles are scaled based on some variable that would not otherwise be plotted. There are cases where this is useful, for instance, adding sample size to a chart via area scaling. Research indicates that people are able to decode area to perceive order but magnitude is difficult accurately translate from area.

**Fig. 5.2** Two views of monthly sales for four departments. The stacked bar chart obfuscates much information that the line chart makes clear



**Choose Graphical Elements Judiciously.** As one builds graphics, there are many choices: colors, shading, line types, line widths, plotting characters, axes, tick marks, legends, etc. Make these choices thoughtfully. Color is often used well when colors indicate membership according to a categorical variable. Color is often used poorly when practitioners get bored and add color haphazardly. Axes can be used intelligently to highlight certain observations or the range of the data. Smoothing lines can illustrate trends in bivariate data or mistakenly cover up observations. With a good graphics package, like the R package `ggplot2` which we introduce in Sect. 5.4, every element of a figure can be manipulated. Take advantage of the opportunities.

**Help Your Audience.** Wherever possible, make adjustments to your figure that helps your audience better understand the data. A choice had to be made in the lower panel of Fig. 5.1. The default behavior sorts the observations by name. This would be the desired order if our goal was to allow the reader to quickly find a given observation in the list and look up the value.<sup>3</sup> But if this is the goal, a table might be a better choice. By sorting the observations by value we can instantly see the minimum and maximum and the observations that define the break between the two groups. By manipulating the data using the R function `reorder`, we help the audience in several ways. Another example of this principle would be to highlight important observations by labeling them.

**Limit Your Scope.** Most interesting projects generate a profusion of data and, as our tools to visualize that data grow, so does the temptation to try to tell the entirety of a story with a single graphic. A well-written paragraph has a topic sentence and a unifying idea. Applying this concept to your graphics requires discipline and attention to detail. There will be times when you are tempted to add additional elements to a chart that already tells the story. Take care that the new elements do not detract from the central message. The classic example of overreaching is a line chart with two different axes for the lines. If you find yourself building such a chart, you are unlikely to be telling the story with clarity and power.

Armed with these general principles, we are now well positioned in the next section to delve into the types of data that are likely to be encountered. We'll describe the elements that make useful visualizations for univariate, bivariate, and multivariate data. In the subsequent section we learn how to produce these graphics in R.

## 5.3 Making Good Choices

Many of our visualization tasks will be defined by the number and type of variables that we will be plotting. The first and most important distinction is between quantitative and qualitative data. A qualitative, or categorical, variable is coded as a `factor` in R. In terms of the number of variables to plot, there are relatively clear approaches when we have one or two variables. Once we move beyond that, there are some general principles in effect but creativity plays a larger and larger role. One must be mindful of the edict to not try to do too much with one chart.

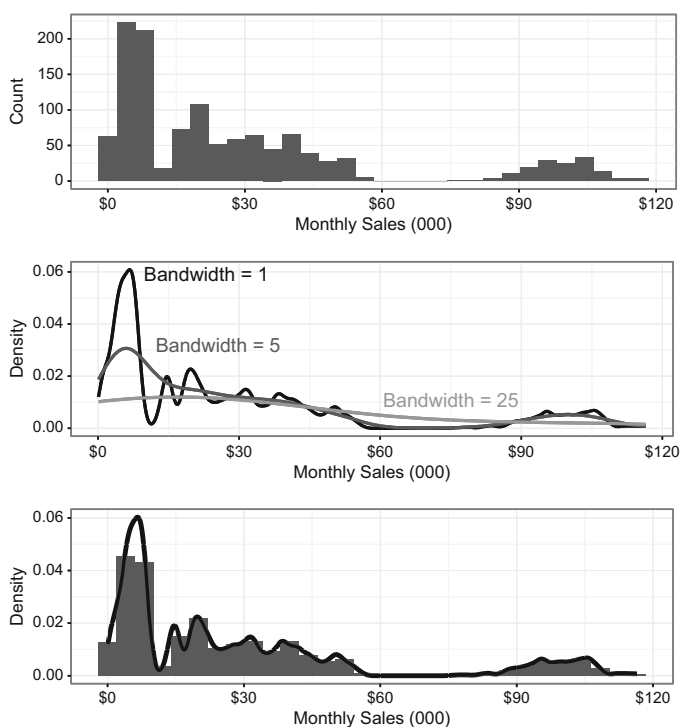
---

<sup>3</sup> When ordering is a problem, it is often referred to as the “Alabama First!” problem, given how often Alabama ends up at the top of lists that are thoughtlessly put, or left, in alphabetical order. Arrange your lists, like your factors, in an order that makes sense.

### 5.3.1 Univariate Data

Our goals with univariate quantitative data typically are to understand the distribution of the data. We typically begin describing the center and spread of the distribution and identifying unusual observations, often called *outliers*. More in-depth analyses describe the shape of the distribution, a task that provides more information and requires more effort. The classical starting point for investigating shape is either the histogram (described at length in Chap. 3, Sect. 3.4.2), the empirical density function, or one of several variations on the boxplot. In this first section we work with a grocery-store data set, specifically, sales and items by month and department for 6 years. We begin by looking at the distribution of sales, in units of thousands of dollars.

**Fig. 5.3** Three different ways of looking at monthly sales by department in the grocery store data: a histogram, several empirical densities illustrating the variations possible from the bandwidth parameter, and a density superimposed on a histogram



The top panel of Fig. 5.3 shows a histogram depicting the numbers of observations falling in each interval, or *bin*, by the height of a vertical bars. We see sales by month by grocery department with bar height representing the count of observations that fall into that “bin”, as the intervals on the  $x$ -axis are referred to.

The second panel illustrates a more powerful and complex way to visualize the distribution of a quantitative variable—empirical density functions. An empirical density function is a data-driven estimate of the probability distribution from which the data originated. More simply, they are smooth histograms. The formula that generates the empirical density function

$$\hat{f}_b(x) = \frac{1}{n \cdot b} \sum_{i=1}^n K\left(\frac{x - x_i}{b}\right), \quad (5.1)$$

where  $n$  is the number of observations,  $b$  is the bandwidth, and  $K$  is the kernel.<sup>4</sup> The bandwidth is usually used to control the smoothness of the function by determining the influence of individual differences  $x - x_i$  on the function. Setting the bandwidth to a small value, say 1, results in a very bumpy distribution. Setting it to a large value arguably makes the distribution overly smooth and removes some of the secondary modes visible with  $b = 1$ . Much like the bin width choice for histograms, some trial and error may be necessary to capture the interesting features of the distribution. The default bandwidth in R which is chosen by `bw.nrd0` and remains the default for historical reasons. The most recommended choice is `bw.SJ`, based on the method of Sheather and Jones [54]. We show how to set the bandwidth in Sect. 5.4.2.

The middle panel of Fig. 5.3 shows three different bandwidths. The bottom panel shows the histogram and density on the same plot, which necessitated changing the units on the histogram to relative frequency. Note that all depictions show some department-months with low sales, a larger group from \$20K to \$55K, and a group at \$100K.

Another useful way to display distributions uses boxplots or violin plots. These plots, illustrated in Fig. 5.4, are more compact displays of the distribution of monthly sales. Boxplots, first developed by John Tukey in the 1970s, summarize the distribution using five numbers. The box is defined by the first and third quartiles,  $Q_1$  and  $Q_3$ , and is split by the median. The interquartile range is  $IQR = Q_3 - Q_1$ . The IQR is a useful nonparametric measure of spread. The whiskers, the length of which can be customized in R, are by default set to  $Q_1 - 1.5 \times IQR$  and  $Q_3 + 1.5 \times IQR$ . Points beyond the whiskers are plotted individually and identified as outliers.

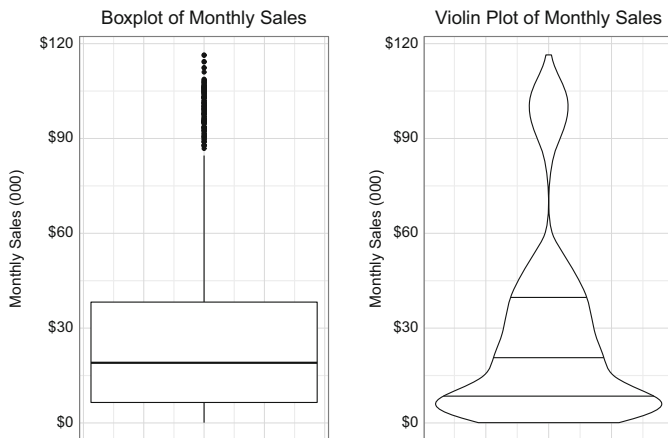
The definition of an outlier may seem arbitrary, but it captures data features in a predictable way if the data are normally distributed. Given a normal distribution, the 25th percentile of the data is  $x_{25} = \mu - .67 \times \sigma$  and the IQR

---

<sup>4</sup> Kernels are an interesting side area of statistics and we will encounter them later in the chapter when we discuss `loess` smoothers. In order for a function to be a kernel, it must integrate to 1 and be symmetric about 0. The kernel is used to average the points in a neighborhood of a given value  $x$ . A simple average corresponds to a uniform kernel (all points get the same weight). Most high-performing kernels use weights that diminish to 0 as you move further from a given  $x$ . The Epanechnikov kernel, which drops off with the square of distance and goes to zero outside a neighborhood, can be shown to be optimal with respect to mean square error. Most practitioners use Gaussian kernels, the default in R.



**Fig. 5.4** Two different ways of visualizing the distribution of monthly sales numbers, the boxplot and the violin plot



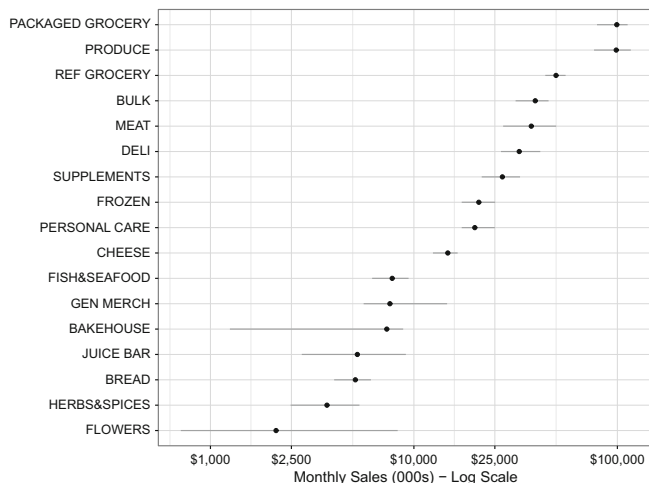
has length  $1.34 \times \sigma$ . Thus,  $1.5 \times \text{IQR} \approx 2\sigma$ . The edge of the whiskers is placed at  $\pm 2.67\sigma$  units from the mean so we would expect less than .8% of the data to fall outside the whiskers. Note that in our example, many points are plotted individually, indicating that the normal distribution is not an appropriate approximation of the monthly sales distribution.

The violin plot, by contrast, makes use of many more features of the data and can be seen as a boxplot replacement. The plot shows a smooth representation of the distribution turned on its side and gives us a more detailed visualization of the distribution. The width of the figure reflects the density of observations. With the addition of the horizontal lines at the first, second, and third quartiles, there is no loss of information compared to the boxplot. Moreover, the most interesting feature of this distribution, the lack of department-months with sales between \$60K and \$80K, is obscured by the boxplot but easy to see with the violin plot.

When we have univariate data that is labeled, we have already seen one of the best visualizations: dotcharts. These charts allow us to clearly depict single values and convey the uncertainty around those estimates (where appropriate). In Fig. 5.1 we saw how a dotchart was superior to a pie chart. In Fig. 5.5 we see another example of a dotchart, this time displaying summary statistics. This figure shows monthly spend by department with bars representing the range of values. Note that we have avoided the “Alabama First!” problem by sorting the departments from highest spend to lowest. In this depiction we can clearly see the two largest departments (produce and packaged grocery), the range of mid-sized departments (refrigerated grocery down to cheese), and the smaller departments. We have translated the  $x$ -axis variable to the  $\log_{10}$  scale. This choice gives us a better view of the monthly sales for the small departments, but can make interpretation a bit trickier for the gray bars, which represent the range. At a first glance, it appears that bakehouse and flowers have by far the widest range in monthly sales. This is

true as a percentage; bakehouse ranges over almost an order of magnitude, from a minimum of \$1000 to almost \$10,000. By contrast, produce has a range of approximately \$40,000 and a sample mean of nearly \$100,000. An important note: the  $x$ -axis displays the values of monthly sales after transforming the values to the  $\log_{10}$  scale. The labels, however, show the original, untransformed monthly sales in thousands of dollars. By retaining the labels in dollars, the reader is better able to interpret the variable. The best practice recommended by Cleveland is to place this axis at the bottom of the graphic and the corresponding log-based scale at the top. Unfortunately, `ggplot2` makes this difficult to do, so our labels show five values on the scale of thousands of dollars.

**Fig. 5.5** A dotchart of spend by month by department, with bars indicating the range of the data. Monthly sales have been transformed to the  $\log_{10}$  scale



The final type of data we might wish to visualize is univariate categorical data. If the number of categories is small, then a simple bar chart is an excellent choice to compare the proportions in each category. If the number of categories is very large, then summarizing with a dotchart as in Fig. 5.5 is often the best practice.

### 5.3.2 Bivariate and Multivariate Data

Bivariate data is usually straightforward to visualize. There are really three main possibilities: two categorical variables, a quantitative and categorical variable, and two quantitative variables.

Data consisting of all categorical variables typically are summarized by contingency tables and there is no reason to deviate from this practice. A contingency table is a cross-tabulation of observations according to the values of

the categorical variables. For instance, if there are  $a$  levels of variable  $A$  and  $b$  levels of variable  $B$ , then the table contains  $a \times b$  cells and the content of a cell is the number of observations with a particular combination of levels. Table 5.1 is a contingency table cross-classifying receipts according to customer segment and department. We can see general trends by examining the table but specific and fine differences are not immediately discernible.

**Table 5.1** Number of receipts cross-classified by department and the three largest customer segments, light, secondary, and primary

Department	Customer segment		
	Light	Secondary	Primary
Supplements	439	55,657	90,017
Cheese	859	96,987	147,679
Frozen	647	97,808	152,940
Meat	653	107,350	149,251
Deli	5830	138,722	155,086
Bulk	2713	144,862	178,520
Refrigerated grocery	3491	194,758	197,463
Produce	3971	211,226	199,978
Packaged grocery	6980	223,815	200,737

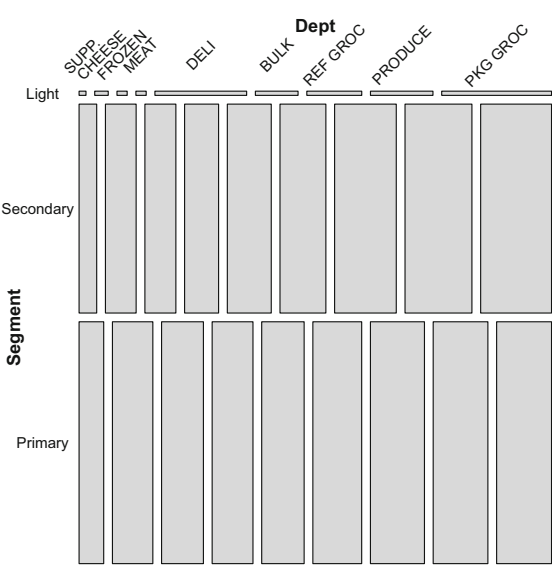
There exists, however, a useful visualization that rapidly conveys the relationships between categorical variables, the *mosaic plot*. Figure 5.6 shows a mosaic plot in which the area of the tiles represents the relative number of observations that fall into each cell of the contingency table. We can see, for instance, that there are many more primary shoppers than light shoppers. Packaged grocery is over-represented among light shoppers, whereas primary shoppers make up larger portions of the less popular departments.

This figure allows us to quickly absorb some of the features of the contingency table:

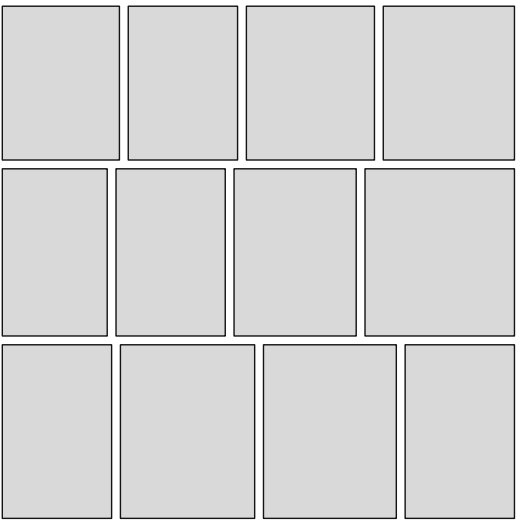
- Primary and secondary shoppers make up the majority of the observations.
- The four largest departments—produce, packaged grocery, refrigerated grocery, and bulk—represent about half the activity of the primary shoppers. Secondary shoppers used those departments to greater extent than primary shoppers and light shoppers used those departments to an even greater extent.
- Supplements represents a small fraction of the purchases of all segments.
- Primary shoppers tend to shop more of the store’s departments, generally.

A virtue of the mosaic plot is that it allows estimation of the strength of the relationship between the categorical variables. A downside is that statistical features, such as confidence intervals for the difference in sizes, cannot be displayed. Moreover, the display becomes unwieldy beyond two dimensions. Figure 5.7 shows a mosaic plot built from data generated from two independent random variables with discrete uniform distributions. The plot shows no evidence of association. Conditioning on one variable shows approximately equal-sized rectangles as you travel across a row or down a column.

**Fig. 5.6** A mosaic plot showing the relationship between customer segments and departments shopped



**Fig. 5.7** A mosaic plot showing no relationship between two categorical variables. Reference plots like this are good to keep in mind when looking at mosaic plots



Most data that we need to visualize is not categorical, however. In the case of a quantitative variable and a categorical variable, we have already seen several good methods of showcasing relationships. Figure 5.5 shows several numeric results (the minimum, mean, and maximum spend by month) split by a categorical variable (the grocery store department). In Fig. 5.8 we see a great deal more information in a similar format. The previous chart showed spend at the department level. This chart shows spend at the individual shopper level for a sample of 10,000 shopper-months.

**Fig. 5.8** A second example of a dotchart showing spending by department at the individual shopper level for 10,000 shoppers

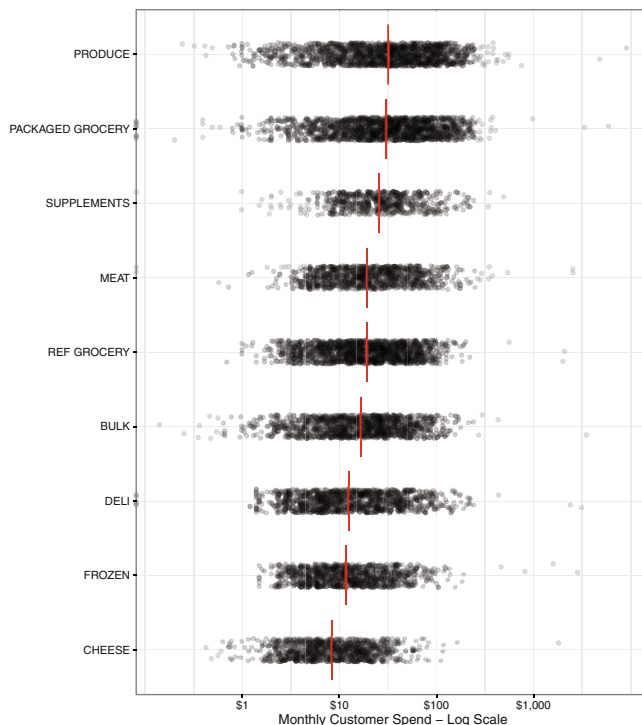
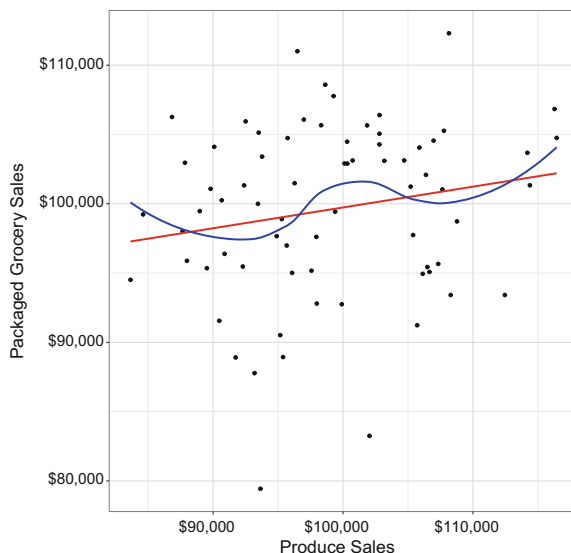


Figure 5.8 displays spend on the  $\log_{10}$  scale by individuals in a month across nine departments in the grocery store. The departments are ordered by the median monthly spend, shown as a vertical red line. The horizontal axis corresponds to a log scale to allow additional detail to be seen at the lower end of the scale. It also appears that the  $\log_{10}$  scale reveals a much more symmetrical distribution—it is not uncommon for retail data to be approximately log-normal in distribution.

Several techniques we have not yet seen are illustrated in Fig. 5.8. The  $y$ -axis shows the levels of department. The values plotted at each level have been jittered so that more of them can be seen. Jittering adds a small random value, say between  $-\varepsilon$  and  $\varepsilon$ , to the vertical coordinate of each plotted pair. Without the jittering all points would collapse onto their nearest horizontal grid line. We have also used transparency, set at the `ggplot2` value  $\alpha = .1$ . The interpretation of this value is that no fewer than  $1/\alpha$  points plotted in the same location will appear completely opaque. Supplements are evidently shopped much less than, say, bulk, but the median spend is higher, presumably because each item is more expensive in this department.

With bivariate numerical data, the natural plotting technique is the scatterplot. Figure 5.9 illustrates this technique.

**Fig. 5.9** A basic scatterplot, showing monthly sales for our two largest departments across 6 years. A linear regression line and loess smoother have been added to the plot to aid in interpretation of the relationship



Scatterplots are straightforward and adhering to the basic tenets laid out in Sect. 5.2 will keep a practitioner on the right track. The data speak for themselves, in this case illustrating the variability between produce and packaged grocery spend when viewed by month. Overall, the department sales move together, with a considerable amount of variation in the relationship.

When illustrating a relationship, one should strongly consider adding a fitted line to the graph to aid the viewer in decoding the relationship. In Fig. 5.9 we add two: a line built by linear regression and a curve created by locally weighted regression. The former needs little elaboration here, as linear regression is covered in Chap. 6. Locally-weighted regression, a technique introduced by Cleveland in 1979 and refined in 1988 [13], is a powerful technique for uncovering the relationship between two numerical variables. There are two competing terms, lowess and loess, and it seems the latter has become preeminent. They are closely related, both setting up a weighted neighborhood around an  $x$  value and fitting a regression line primarily influenced by points in vicinity of  $x$ . Much like our kernel smoothers discussed above, a bandwidth parameter,  $\alpha$ , is specified. In each neighborhood a polynomial of degree  $d$  will be fit,<sup>5</sup> and we require a choice of  $\alpha \in [\frac{d+1}{n}, 1]$ . The default settings in R are  $\alpha = .75$  and  $d = 2$ . For each subset of size  $n\alpha$ , a polynomial is fit. This fit is based on weighting the points and the typical weight function is the tricube weight. Let us assume that we have an observation  $x_i$ , a neighborhood around  $x_i$  denoted by  $N(x_i)$ , the width of which is  $r_i$ . Then the weight function is

<sup>5</sup> In practice  $d$  is almost always 1 or 2.

$$w_i(x) = \begin{cases} \left(1 - \left|\frac{x-x_i}{r_i}\right|^3\right)^3, & x \in N(x_i), \\ 0, & x \notin N(x_i). \end{cases}$$

The tricube kernel is nearly as effective as the Epanechnikov kernel. In practice, the choice of kernels is not nearly as important as the bandwidth choice.

We now turn our attention to multivariate data. Multivariate data is defined as data comprising at least three covariates per observation and, as mentioned earlier, most solutions require a thoughtful process. We can apply the general principles mentioned earlier and enjoy the profusion of options that are available. If there is one principal that stands above all others it is to avoid doing too much. This temptation is nearly irresistible. One may violate the rule profitably when a graphic will do the heavy lifting for several pages of text or when the graphic can be displayed on a slide for several minutes of explanation. This is not the norm and data scientists would be wise to split their story into multiple graphics if the audience is pressed for time.

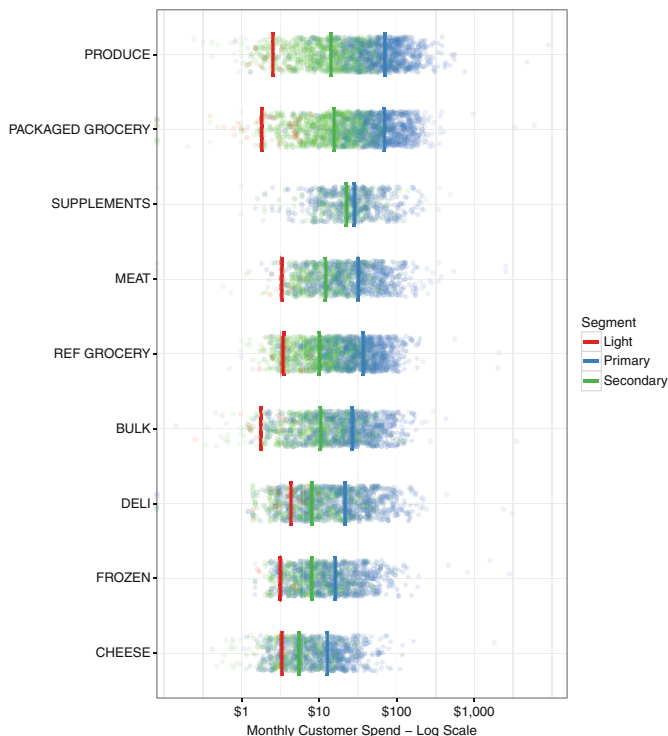
Figure 5.10 is a reprisal of Fig. 5.8, made multivariate with the addition of the customer segmentation of the shopping data seen in Chap. 10. A number of interesting features emerge, notably the separation between secondary and primary shoppers for all departments except for supplements. The addition of median lines for each segment aids comprehension, particularly for the light shoppers, who make up small fraction of the shoppers. Interestingly, we can now see the compression of the deli spend—this department is one of the most popular with light shoppers.

Figure 5.10 is, essentially, a two-dimensional data visualization with a third dimension (segment), layered on top. Using two-dimensional plots with additional information to encode other variables is a common technique. Another variation, illustrated by Fig. 5.11 splits a two-dimensional plot into small multiples of pairs according to a third variable. The term “small multiples”, coined by Tufte, captures the idea that readers, once they have been oriented to an individual plot, can quickly discern similarities and differences across a family of plots.

Each small panel, known as a “facet”, is a scatterplot of spend versus items within a given department. A loess smoother is added to each panel. The advantages of faceting versus simply repeating scatterplots are several-fold: parsimonious code, an efficient use of space on the layout, the ability to order the facets in a sensible way, and common axes that facilitate comparison. With this treatment we can quickly identify interesting patterns within the data:

- The large volume of sales in produce, packaged grocery, and refrigerated grocery stand out relative to the other departments.
- Steeper curves indicate departments with cheaper items (produce, deli) while flatter curves show the expensive items (supplements, meat).
- Certain departments do not have large spends (cheese, frozen, and bulk).

**Fig. 5.10** An example of multi-variate data. The spend-department dotchart is now colored based on the segments of the shoppers (either primary, secondary, or light). The medians are shown for each segment as a color-coded line



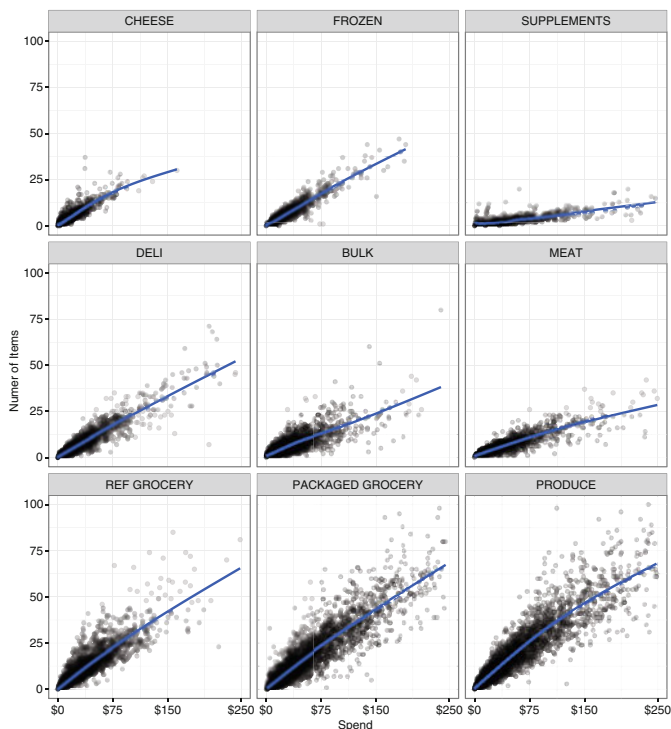
This section has served as a travelogue without a map. We have seen the destinations, but now we must learn how to get there. Building these plots requires a surprisingly small amount of code using the package that we will learn, `ggplot2`, partially because a great deal of complexity is hidden from us. Graphical software is a leaky abstraction, unfortunately, and so we will illustrate how to build up these charts from base elements.

## 5.4 Harnessing the Machine

The plots in this chapter were made with the R package `ggplot2`. This package is our strong recommendation for your personal data visualization tasks. There are three widely-used graphics packages associated with R: `base`, `lattice` [53], and `ggplot`. The `base` package is useful for quick plotting and for learning the basic techniques. It is possible to control many aspects of the plot with the `base` package but constructing publication-quality graphics is not easy. The second package, `lattice`, is based very closely on the ideas of Cleveland, but is not developed on the framework of a formal model. The `lattice` framework has limited its extensibility. The package `ggplot` is more flexible and easier to use.



**Fig. 5.11** A scatterplot, faceted by department, of spend versus items, with a loess smoother added to each panel



The package name, `ggplot2`, reflects both the version number, two, and the heritage, Wilkinson’s Grammar of Graphics [69]. To quote Hadley Wickham, the principal author of `ggplot2`, “Wilkinson created the grammar of graphics to describe the deep features that underlie all statistical graphics.” [66] Wickham’s books covers the grammar in some detail. Our treatment draws extensively from Wickham’s texts.

The `ggplot` grammar requires work on the part of the reader. Once it is mastered, however, `ggplot2` provides the user with lots of flexibility and power. We encourage the reader to learn the grammar.

The grammar of graphics consists of a number of components that merge to form the grammar. They are

1. **data** The data to be rendered as a visual. Using `ggplot`, the data must be an R data frame. A R data frame is rectangular arrangement of the data with somewhat more specificity and overhead than a simple matrix. For example, a data frame may contain variables of several types, say, quantitative (numeric in the R lexicon) and categorical (a factor in the R lexicon).
2. **aes** *aes* is shorthand for *aesthetic mapping*. Aesthetic mappings tell `ggplot2` how to translate the data into graphic elements. The aesthetic mapping identifies the variables to be plotted on the *x*- and *y*-axes.

More may be added to the aesthetic mapping. For example, a categorical variable that determines the color of the points or lines is identified in the aesthetic mapping.

3. **geoms** *geoms* is shorthand for *geometric objects*. These are the elements that visually portray the data. Examples are points, lines, line segments, error bars, and polygons. All geometric object specifications take a form such as `geom_points()`.
4. **stats** These are statistical transformations that are used to reduce and summarize the data. The smoothers we saw above are examples of statistical transformations since the data was reduced in some manner to produce the smooth lines.
5. **scales** Scales provide the mapping between the raw data and the figure. Scales are also used to create legends and specialized axes.
6. **coord** The coordinate system takes us from the data to the plane of the graphic. This component is used, notably, to change axis scales from the original units of a variable to different units, say, logarithms.
7. **facet** As we saw in Fig. 5.11, facetting divides the graphic into sub-graphics according to a categorical variable. This graphical component defines how the facets are arranged.

The process of building a visualization as a series of *layers* in `ggplot` is a remarkable improvement on the traditional process. Figures in `ggplot` are built by first creating a base consisting of a rudimentary plot and then adding more information and data to the base in the form of layers. Each layer may contain specific information from the attribute list above. Usually, layers inherit most of the attribute values from the initially created plot. The specification of the layer may be very simple since we only need to change a few items in each layer. Consequently, layers make the task of building a complex graphic much easier. We tend to build graphics one layer at a time. We can see the effect of each layer on the graphic and more easily correct coding errors. A highly readable tutorial on the subject is available at <https://rpubs.com/hadley/ggplot2-layers>.

Our intent in this chapter is not to systematically review `ggplot2` techniques but to provide a general understanding of how visualizations are constructed in `ggplot2`. In the remainder of the chapter, we explain how the graphics discussed above were built. To go further, utilize internet resources and, in particular, Stack Overflow (<https://www.stackoverflow.com>). If you can't solve a problem or remember an instruction, then search the internet for information. Adding *ggplot* to the search string makes one much more likely to find results related to R and to the `ggplot2` plotting package in particular. Using the vocabulary of the package in the search string is important to efficient searching.

We now turn to the code that created this chapter's figures. We first start by reading in the data and loading the necessary libraries.

```
library(ggplot2)
library(scales)
library(RColorBrewer)
```

The first two libraries, `ggplot2` and `scales` are directly related to plotting. The final one, `RColorBrewer`, is based on the pioneering work on color and perception of Brewer et al. [7], is indispensable and we recommend it highly.

### 5.4.1 Building Fig. 5.2

Figure 5.2 is built from a summary table named `month.summary`. This data frame holds the sum of monthly sales for four departments at the co-op. The first five rows are displayed in Table 5.2.

**Table 5.2** The first few rows of the data frame `month.summary`

Row	Month	Department	Sales
1	6	Packaged grocery	90,516.72
2	6	Produce	103,650.8
3	6	Bulk	37,105.20
4	6	Meat	35,221.97
5	5	Packaged grocery	102,178.30

The code follows.

```
month.summary <- read.delim("../data/month_summary.txt")

ggplot(data=month.summary,
       aes(x=factor(month), y=sales/1000, group=Dept, col=Dept) )
+ scale_color_brewer(palette="Set1")
+ geom_point()
+ geom_line()
+ theme_bw()
+ scale_y_continuous(label=dollar)
+ ylab("Avg Annual Sales in Month (000s)")
+ xlab("Month")
+ theme(legend.key.size=unit(0.5, "cm") )
```

After reading in the data, we build the plot in `ggplot2`, layer by layer.

1. **data** We use the data shown in Table 5.2.
2. **aes** We assign month to the  $x$ -axis and sales-divided-by-one-thousand to the  $y$ -axis. Department is used as a grouping variable. The consequence of

setting `group=Dept` is that points belonging to the same department will be joined by lines. Departments will be identified by color. When building graphics that use grouping and color, it is common to forget to include the grouping variable in the aesthetic. As an exercise, we encourage the reader to build the plot *without* declaring the grouping variable in `aes`. The appearance of the resulting plot, with the telltale diagonal lines, are the mark of a missing grouping variable.

3. **geoms** There are two geometric objects used in the figure: points and lines. They are added as separate layers and inherit the grouping and color from the aesthetic.
4. **scales** We use a continuous scale and, for the axis labeling, take advantage of the library `scales` to reduce the effort of the reader to understand what has been plotted. The units are dollars so we pass the keyword `dollars` into the `label` argument. This parameter will ensure that the axis labels are formatted with dollar signs, commas, and cents. Other options include `percent` and `comma`, all of which promote readability. We also use the colors provided by `RColorBrewer` for the same reason.
5. **theme** Although not part of our original list, themes help you adjust the appearance of your plots. We invoke the simple and clean `theme_bw`, label the axes, and shrink boxes in the legend a bit. These features are almost always included in our graphics.

### 5.4.2 Building Fig. 5.3

Our next code segment builds our histogram and empirical density function.

```
# Read in grocery store summary data.
working.dir <- "../data/"
gd <- read.delim(paste0(working.dir, "grocery_data.txt"))
ggplot(gd, aes(x = ownerSales/1000))
  + geom_histogram(aes(y=..density..))
  + geom_line(stat="density", col="gray50", size=1.5,
             bw="SJ") + # Changing the bandwidth algorithm to 'SJ'
  + theme_bw()
  + scale_x_continuous(label=dollar)
  + ylab("Density")
  + xlab("Monthly Sales (000)")
```

The key features used to construct the figure are as follows.

1. **data** Histograms summarize distributions as they are built. We don't have to carry out data reduction before building the figure. It's done by `ggplot` in the construction of the histogram.

2. **aes** Histograms are defined by a single variable which we assign to the  $x$ -axis. `ggplot2` determines the  $y$ -values using our guidance.
3. **geoms** The first geometric object is a histogram. The **aes** argument sets the scale for the  $y$ -axis. Our choice of units for the histogram are proportions. The  $y$ -axis will show the proportion of observations in each interval and it's specified by setting `y=..density..`. The syntax of `..density..` looks odd. The pattern of periods (two at the beginning and two at the end) indicates a statistical computation is necessary. The other commonly chosen option shows the counts of observations in each interval. Entering `?stat_bin` from the console will provide more information.

A second geom, `line` is also specified. The **stat** argument specifies that the line to be drawn is the graph of an empirical density function (Sect. 5.3.1). Several other attributes are specified: the color, the line width (by adjusting `size`), and that the method of computing the bandwidth is to be Sheather and Jones' method [54].

### 5.4.3 Building Fig. 5.4

The code for the violin plot from Fig. 5.4 is next.

```
ggplot(gd, aes(x=1, y=ownerSales/1000))
+ geom_violin(draw_quantiles = c(0.25,0.5,0.75))
+ theme_bw()
+ scale_y_continuous(label=dollar)
+ xlab("")
+ ylab("Monthly Sales (000)")
+ theme(axis.text.x=element_blank(),axis.ticks=element_blank())
+ labs(title="Violin Plot of Monthly Sales")
```

A couple of new techniques emerge.

1. **aes** The  $x$ -axis is set to be a constant. The result is that the violin plot will be drawn in the center of the graphic.
2. **geoms** The `geom_violinplot` accepts an argument specifying that quantiles are to be drawn. We specify that these are to be the first, second, and third *quantiles*.
3. **theme** Since the  $x$ -axis does not have any meaning, we would like to avoid showing ticks marks and labels. The arguments to **theme** eliminate those features. A internet search for something similar to “ggplot blank x axis” will provide details.

### 5.4.4 Building Fig. 5.5

The dotchart in Fig. 5.5 requires a prepared data frame `dept.summary` containing the sample minimum, mean, and maximum sales the 17 departments of the grocery store co-op. The first five rows are shown in Table 5.3.

**Table 5.3** The first five rows of the `dept.summary` data.frame

Row	Department	Sample		
		Minimum	Mean	Maximum
1	Packaged grocery	79,434.24	99,348.55	112,303.60
2	Produce	76,799.20	98,711.50	116,442.92
3	Bulk	31,610.05	39,482.60	45,908.04
4	Refrigerated grocery	44,239.77	49,940.93	55,656.07
5	Cheese	12,407.99	14,674.32	16,404.14

After reading the summary data into a data frame, it is re-ordered by the sample mean.

```
dept.summary <- read.delim("../data/dept_summary.txt")
dept.summary$dept <- reorder(dept.summary$dept, dept.summary$mean.val)

ggplot(dept.summary, aes(x=mean.val, y=dept))
  + theme_bw()
  + geom_errorbarh(aes(x=mean.val, xmax=max.val, xmin=min.val, y=dept),
    height=0, color="gray60")
  + geom_point(col="black")
  + ylab("")
  + xlab("Monthly Sales (000s) - Log Scale")
  + scale_x_continuous(label=dollar,trans="log10",
    breaks=c(1000,2500,10000,25000,100000))
```

1. **aes** When a factor is used as an  $x$ - or  $y$ -variable in an aesthetic, the factor levels are mapped to sequential integers, say  $1, 2, \dots, g$ , where  $g$  is the number of levels. Knowing that the levels occur at integer positions on the  $x$ -axis will be important when we add features like jittering.
2. **geoms** Another new `geom`, `geom_errorbarh`, is used. The “h” stands for horizontal—the vertical variety needs no suffix. The aesthetic mapping for `geom_errorbarh` requires us to supply  $x$ - or  $y$ -variables and the starting and ending positions for the error bars. The parameter `height` specifies the width of the whiskers on the bars. We’ve suppressed the whiskers as they add no information. Note that points are desired at the median and receive their own `geom`.
3. **scale\_x\_continuous** We repeat our label trick and instruct `ggplot` to transform the variable plotted on  $x$ -axis according to the transformation  $x \rightarrow \log_{10}(x)$ . Also, we set the label positions using the `breaks` argument.

### 5.4.5 Building Fig. 5.8

We build the plots for Figs. 5.8 and 5.10 in the next series of code segments. We begin using `sample.int` to draw a random sample of manageable size from the shopper data set. Setting the seed of the random number generator with the instruction `set.seed` ensures that our results are the same on repeated runs. We also order our departments based on the median total sales by department.

```
set.seed(3939394)
# Read in grocery store detail data.
gdd <- read.delim(paste0(working.dir,"shopper_dept_segment.txt"))
this.data <- gdd[sample.int(nrow(gdd),size=10000,replace=F),]
this.data$DepartmentName <- reorder(this.data$DepartmentName,
                                   this.data$TotalSales,
                                   FUN=median)
```

Drawing the vertical lines that depict the medians requires those values to be calculated and it is clearest if we assign them to a data frame of their own. We call that data frame `medians`.

```
medians <- aggregate(this.data$TotalSales,
                     list(this.data$DepartmentName),
                     FUN=median)
names(medians) <- c("DepartmentName", "TotalSales")
medians$y.val <- as.numeric(medians$DepartmentName)
```

We are now ready to build Fig. 5.8.

```
ggplot(this.data, aes(x = TotalSales, y = DepartmentName))
+ theme_bw()
+ geom_point(position=position_jitter(h = .4),
             alpha=0.1)
+ ylab("")
+ xlab("Monthly Customer Spend - Log Scale")
+ scale_x_continuous(label=dollar, trans="log10",
                     breaks=c(1,10,100,1000))
+ geom_segment(data = medians,
               aes(x = TotalSales, xend = TotalSales,
                   y = y.val -.4, yend = y.val+ .4),
               col="red")
```

1. **data** Two data sets are used to construct the figure. The principal data set is the sample of individual shopper grocery data. A second data set called `medians` is used for the vertical lines and is passed directly into the necessary `geom_segment`.

2. **aes** This plot uses an aesthetic that is similar to what we have seen before, with the categorical variable `departmentName` assigned to the  $y$ -axis. Individual shopper sales are assigned to the  $x$ -axis. There is a separate set of aesthetic mappings used with `geom_segment`.
3. **geoms** There are two **geoms** used in this plot, one for the points and one for the line segments. The `geom_point` has been used before, but now we show how to add jittering to a plot. Recall that jittering perturbs points so that there is less over-plotting and individual points are easier to see. The command, `position=position_jitter(h=.4)` instructs `ggplot2` to place the points at the original  $x$ -axis coordinate and to perturb the  $y$ -axis coordinate. The argument  $h = .4$  controls the magnitude of the random perturbations.  
The other **geom** layer, `geom_segment` uses the second data set. We specify the starting and ending line segment coordinates on the  $x$ - and  $y$ -axes.
4. **scale** Once again we use the *dollar* labeling from the library `scales`. As we have been doing in other plots, we instruct `ggplot` to transform the variable plotted on  $x$ -axis according to the mapping  $x \rightarrow \log_{10}(x)$  and specify where we would like the labels to appear.

### 5.4.6 Building Fig. 5.10

The code used to construct Fig. 5.10 is a relatively simple modification of the code used to construct Fig. 5.8. The key difference is that we wish to color the points according to segment, as well as have separate median lines for the segments. First, we must calculate the medians for each department and segment.

```
medians <- aggregate(this.data$TotalSales,
                     list(Segment=this.data$Segment,
                          DepartmentName=this.data$DepartmentName),
                     FUN=median)
names(medians)[3] <- "TotalSales"
jit.val <- 0.6 # More jittering with segments
```

The function `aggregate` is convenient for summarizing one variable based on one or more other variables. We rename the column that holds the calculated medians to match the aesthetic we'll use in our plot.

As we shall see, it is quick to add a grouping variable that allows coloring of the points and the technique is analogous to what we did in Fig. 5.2. There we used grouping and color aesthetics to ensure that the department sales by month were joined by a line and common color. Here the variable `Segment` will fill that role.



```
ggplot(this.data,
      aes(x=TotalSales, y=DepartmentName, group=Segment, col=Segment))
+ theme_bw()
+ geom_point(position=position_jitter(h=jit.val),
             alpha=0.1)
+ ylab("")
+ scale_color_brewer(palette = "Set1")
+ xlab("Monthly Customer Spend - Log Scale")
+ scale_x_continuous(label=dollar,trans="log10", breaks=c(1,10,100,1000))
+ geom_segment(data = medians,
              aes(x = TotalSales,xend=TotalSales, group=Segment, col=Segment,
                  y = y.val-0.5*jit.val, yend=y.val+0.5*jit.val), size=1.25)
```

What are the differences between this code and the previous code? Just the addition of the argument `group=Segment` and `col=Segment` in the `ggplot` aesthetic and `geom_segment` aesthetics, respectively. Other aspects of our plots remain the same.

### 5.4.7 Building Fig. 5.11

Our final example builds Fig. 5.11. The code is reproduced here, though there are few features we have not seen before.

```
set.seed(3939394)
this.data <- gdd[sample.int(nrow(gdd),size=25000,replace=F),]
this.data$DepartmentName <- reorder(this.data$DepartmentName,
                                   this.data$TotalSales,
                                   FUN=sum)
ggplot(this.data, aes( x = TotalSales, y = Items))
+ geom_point(alpha=0.1)
+ stat_smooth(se=F)
+ facet_wrap(~DepartmentName)
+ scale_x_continuous(label=dollar, limits=c(0,250),
                    breaks=c(0,75,150,250))
+ scale_y_continuous(limits=c(0,100),label=comma)
+ theme_bw()
+ xlab("Spend")
+ ylab("Number of Items")
```

The first new feature is a statistic that stands on its own (as opposed to being hidden inside `geom_line` as we saw above). The function, `stat_smooth`, is quite useful, adding a line to a plot according to the specified method. In the code that generated Fig. 5.9, we added two smooths in two layers. One layer added the default smoothing method, loess, and the second added a fit-

ted linear regression line using the instruction `stat_smooth(method="lm")`. The default is to include standard error envelopes. We have suppressed the envelopes by passing the argument `se=F`.

The other new feature is `facet_wrap`, the function that controls the creation of the panels. The *wrap* version creates a rectangular array of panels. One panel is created for each level of the factor `DepartmentName`. Only observations associated with a particular level of the factor are graphed in the panel. Rows and columns can be controlled. Multiple variables can be set on the right-hand side of the tilde (`~`) to partition the data by combinations of levels of the variables. The analog of `facet_wrap` is `facet_grid`, which is designed to set up the panels in a tabular array, specifying row factors and column factors using the same tilde operator. One important parameter, not changed in this code, is the `scales` value, which can be set to `free`, `free_x`, `free_y`, and `fixed`. The default is the last one and it forces a single set of axes for every panel. The `free` varieties allow the axes to differ and generates the best visual fit for each panel. The `free` option is often a mistake if the purpose of the panels is to compare and contrast groups because the reader will have to account for different scales in different panels. The best visualization depends on the goals of the data scientist.

A final word of caution and encouragement. Plotting with `ggplot2` imposes an initially steep learning curve. The payoff is well worth the effort—do not underestimate the importance of communicating clearly and efficiently. Visualization is the best way to do it. The key to learning `ggplot` is to start simple and add complexity one layer at a time. Build a simple scatterplot with `geom_point`, then add a layer that identifies groups by color. Add a layer showing smooths. Label the axes clearly and adjust the font sizes. Every successful plot that we created for the first 6 months of working with `ggplot2` was made by incremental steps.

## 5.5 Exercises

**5.1.** Use the data set `grocery_data.txt` containing sales by month and build a dotchart of monthly sales by department. Put department on the *y*-axis and sales (in thousands) on the *x*-axis. Include the following features:

- Sort the departments based on average sales.
- Jitter the points vertically and set `alpha` to help with over-plotting.
- Label the *x*-axis using dollars.
- Suppress the *y*-axis label and use a sensible label for the *x*-axis.

What patterns emerge from the departments? How does the spread of sales vary by department? Why might this be?

**5.2.** Add vertical red lines representing the median monthly sales by department to the figure of Problem 5.1.

**5.3.** Build the middle panel from Fig. 5.3 with the bandwidth parameter set to the values 1, 5, and 25. What is the bandwidth returned by the “SJ” method? How does it compare to the `bw.nrd0` method<sup>6</sup>?

**5.4.** Continuing to use the data set `grocery_data.txt`, build a faceted graphic showing the empirical density of sales per month by department. Allow the  $x$ - and  $y$ -axes to vary within the panels. Use the SheatherJones bandwidth selection algorithm. Compare the densities.

**5.5.** Now we fine-tune the solution to Problem 5.4 to make it publication worthy. Add the following features:

- a. Angle the  $x$ -axis labels to  $45^\circ$ . (Hint: `?theme` and search for `axis.text.x`.)
- b. Use the black and white theme.
- c. Remove the beer & wine segment from the plot since there are very few observations in the segment.
- d. Shrink the font size of the facet titles (called `strip.text`) to make the longer names fit.

**5.6.** Repeat Problem 5.1 with violin plots. Add quartile lines and draw a red vertical line at the overall average monthly sales across all months and all departments. Use the parameter `scale="width"` in the violin plot to avoid the default behavior of the width being proportional to sample size.

Does this view of the data add to any insights beyond what you learned from Problem 5.1?

**5.7.** Build the mosaic plot shown in Fig. 5.6. This plot is *not* made by `ggplot`, but is made with the `mosaic` function in the package `vcd`. Controlling the orientation of the axis labels is tricky; we recommend replacing the department and segment names with abbreviations.

**5.8.** Use the `gdd` data to build the sum of sales by year and month, and segment. Display these in a line chart in chronological order, with each segment being its own line.

**5.9.** Use the `gdd` data to build the sum of sales by year and month and segment. Plot sum of sales against year and month, with each segment shown by a separate line. Include the following features in your plot:

- a. Vertically align the year-month labels.
- b. Plot the  $y$ -axis on a log scale and choose sensible breakpoints. Label the axis with the “dollar” formatting.
- c. Label the  $x$ -axis every 12 months starting at 2010-1.
- d. Add a linear trend line, by segment, to the plot. What information can you infer about the light segment from this graph?

---

<sup>6</sup> The purpose of these last two questions is to learn how to extract the bandwidth information from R directly.