





# Representing Data

When you visualize data, you represent it with a combination of visual cues that are scaled, colored, and positioned according to values. Dark-colored shapes mean something different from light-colored shapes, or dots in the top right of a two-dimensional space mean something different than dots in the bottom left.

Visualization is what happens when you make the jump from raw data to bar graphs, line charts, and dot plots. It's the process that takes you from the grid of photos in Chapter 1, "Understanding Data," to a bar graph over time, as shown in Figure 3-1.

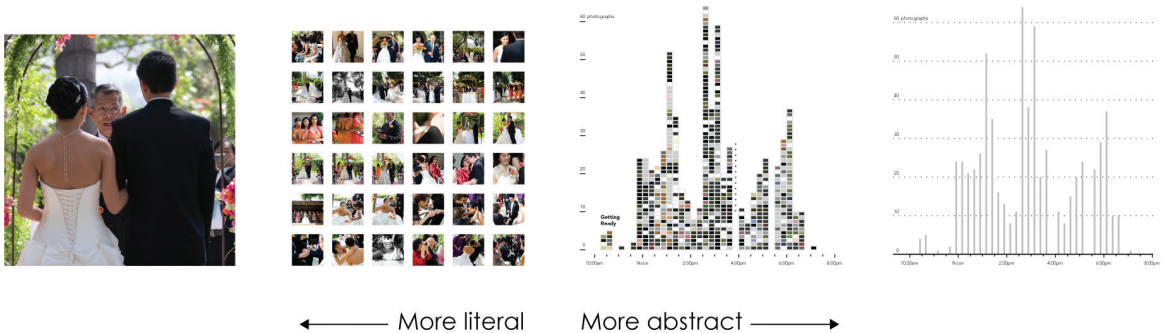


FIGURE 3-1 *Abstraction process*

It's easy to think that this process is instant because software enables you to plug data in, and you get something back instantly, but there are steps and choices in between. What shape should you choose to encode your data? What color is most appropriate for the purpose and message? You can let the computer choose everything for you (it can save time), but there are advantages when you choose. At the least, if you know the elements of visualization and how they can be combined and modified, you know what to tell the computer to do rather than let the computer dictate everything you make.

In many ways, visualization is like cooking. You are the chef, and datasets, geometry, and color are your ingredients. A skilled chef, who knows the process of how to prepare and combine ingredients and plate the cooked food, is likely to prepare a delicious meal. A less skilled cook, who heads to the local freezer section to see what microwave dinners look good, might nuke a less savory meal. Of course, some microwave dinners taste good, but there are a lot that taste bad.



Whereas the person who is only familiar with entering the time and power level on a microwave must either endure poor-tasting meals or stick only to the handful of good ones, people who understand the ingredients and actually know how to cook have fewer limitations. The skilled chef might even transform an average frozen dinner into a gourmet meal.

Likewise, with visualization, when you know how to interpret data and how graphical elements fit and work together, the results often come out better than software defaults.

## VISUALIZATION COMPONENTS

What are the ingredients of visualization? Figure 3-2 shows a breakdown into four components, with data as the driving force behind them: visual cues, coordinate system, scale, and context. Each visualization, regardless of where it is on the spectrum, is built on data and these four components. Sometimes, they are explicitly displayed, and other times they form an invisible framework. The components work together, and your choice with one affects the others.

**Note:** Cartographer Jacques Bertin described a similar breakdown in *Semiology of Graphics*, and statistician Leland Wilkinson later provided a variation in *The Grammar of Graphics*.

### VISUAL CUES

In its most basic form, visualization is simply mapping data to geometry and color. It works because your brain is wired to find patterns, and you can switch back and forth between the visual and the numbers it represents. This is the important bit. You must make sure that the essence of the data isn't lost in that back and forth between visual and the value it represents because if you can't map back to the data, the visualization is just a bunch of shapes.

You must choose the right visual cue, which changes by purpose, and you must use it correctly, which depends on how you perceive the varied shapes, sizes, and shades. Figure 3-3 shows what's available.

#### Position

When you use position as a visual cue, you compare values based on where others are placed in a given space or coordinate system. For example, when you look at a scatterplot, as shown in Figure 3-4, you judge a data point based on its x- and y-coordinate and where it is relative to others.

**FIGURE 3-2** (following page)  
The components of visualization

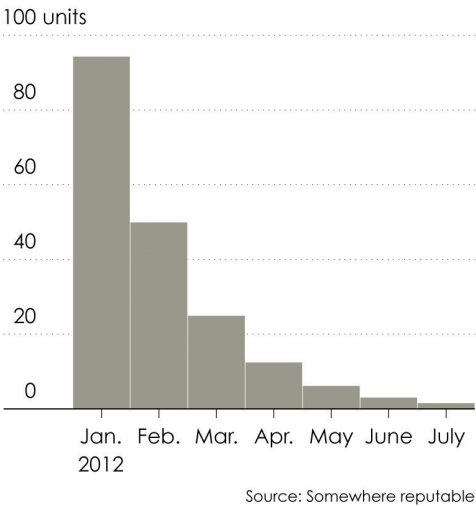


# Working parts

Several pieces work together to make a graph. Sometimes these are explicitly shown in the visualization and other times they form a visual in the background. They all depend on the data.

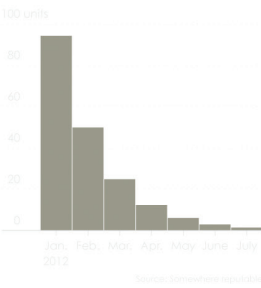
## Title of this Graph

A description of the data or something worth highlighting to set the stage.



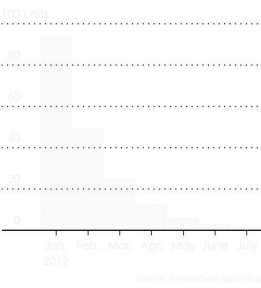
## Title of this Graph

A description of the data or something worth highlighting to set the stage.



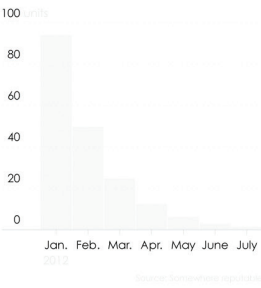
## Title of this Graph

A description of the data or something worth highlighting to set the stage.



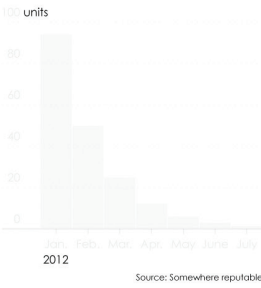
## Title of this Graph

A description of the data or something worth highlighting to set the stage.



## Title of this Graph

A description of the data or something worth highlighting to set the stage.



# Visual Cues

Visualization involves encoding data with shapes, colors, and sizes. Which cues you choose depends on your data and your goals.

# Coordinate System

You map data differently with a scatterplot than you do with a pie chart. It's x- and y-coordinates in one and angles with the other; it's cartesian versus polar.

# Scale

Increments that make sense can increase readability, as well as shift focus.

# Context

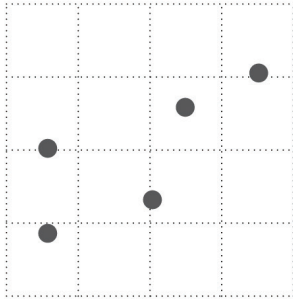
If your audience is unfamiliar with the data, it's your job to clarify what values represent and explain how people should read your visualization.

# Visual cues

When you visualize data, you encode values to shapes, sizes, and colors.

## Position

Where in space the data is



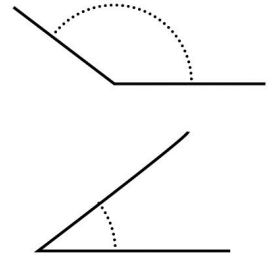
## Length

How long the shapes are



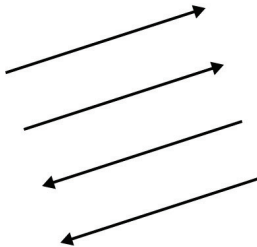
## Angle

Rotation between vectors



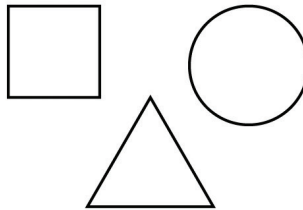
## Direction

Slope of a vector in space



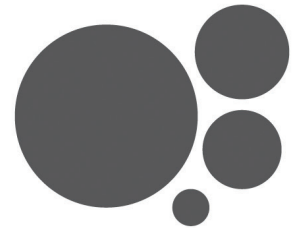
## Shapes

Symbols as categories



## Area

How much 2-D space



## Volume

How much 3-D space



## Color saturation

Intensity of a color hue



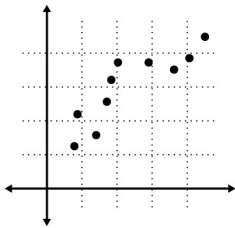
## Color hue

Usually referred to as color

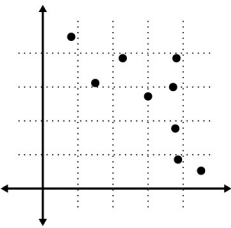


**FIGURE 3-3** Visual cues

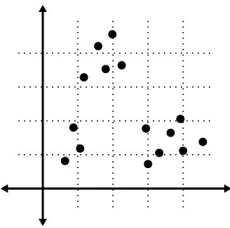
### Upward trend



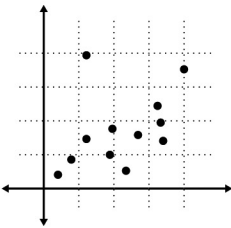
### Downward trend



### Clustering



### Outlier



**FIGURE 3-4** Scatterplots

One of the advantages of using just position is that it tends to take up less space than the other visual cues because you can draw all the data within the x- and y-plane, and you can represent each point with a dot. Unlike other visual cues that use size to compare values, all points in a position-based plot are the same size. In turn, you can spot trends, clusters, and outliers by plotting a lot of data at once.

However, the advantage of using position alone can also be a disadvantage. If you look at a lot of points at once in a scatterplot, it can be a challenge to identify what each point represents. Even in an interactive plot, you still must mouse over or select a point to find out more information, and overlap can cause more problems.

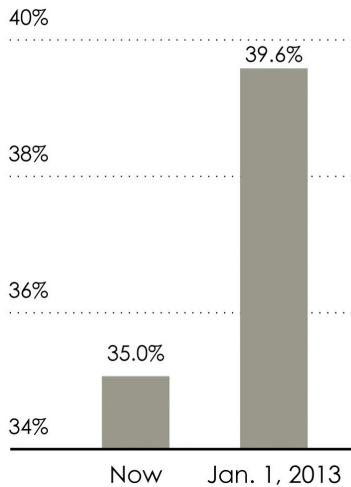
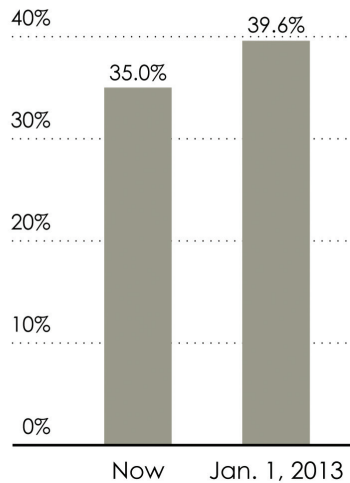
### Length

Length is most commonly used in the context of bar charts. The longer a bar is, the greater the absolute value, and it can work in all directions: horizontal, vertical, or even at different angles on a circle.

How do you judge length visually? You figure out the distance from one end of a shape to the other end, so to compare values based on length, you must see both ends of the lines or bars. Otherwise, you end up with a skewed view of maximums, minimums, and everything in between.

As a simple example, as shown in Figure 3-5, a major news outlet displayed a bar graph on television that compared a tax rate before and after a date.

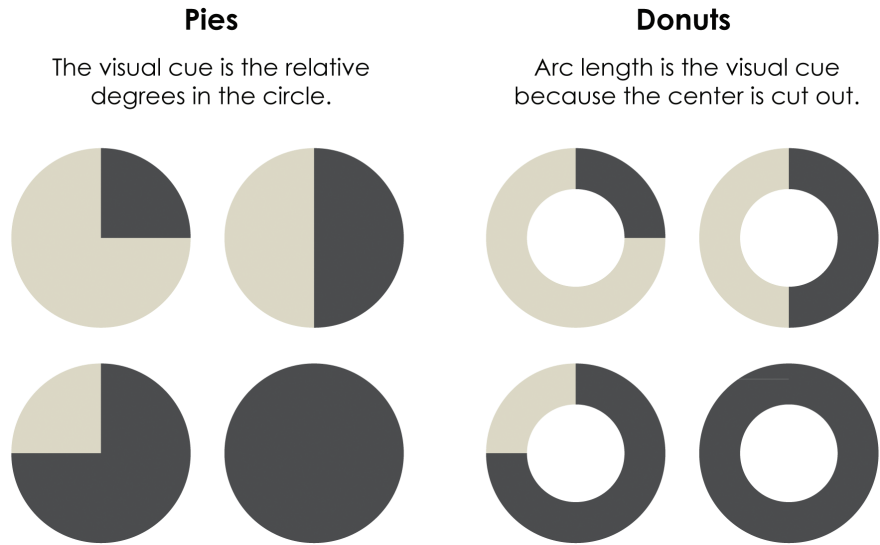


**Axis starting at 34 percent****Axis starting at 0 percent****FIGURE 3-5** *Incorrect bar graph on left and correct one on the right*

The difference between the two values looks like a huge increase—the length of the right bar is about five times the length as the other—because the value axis starts at 34 percent. The chart on the right shows the change when the axis starts at zero, which looks less dramatic. Of course, you can always look at the axis to verify what you see (and you always should), but that defeats the purpose of showing the values with length, and if the chart is shown quickly on television, most people won't notice the misstep.

### Angle

Angles range from zero to 360 degrees on a circle. There's the 90-degree right angle, the obtuse angle that is greater than 90 degrees, and the acute angle that is less than 90 degrees. A straight line is 180 degrees.



**FIGURE 3-6** Mmm, pies and donuts

**Note:** Although the donut chart is often considered the pie chart's close cousin, arc length is the former's visual cue because the center of the circle, which indicates angles, is removed.

For each angle in between zero and 360 degrees, there is an implied opposite angle that completes the rotation, and together those two angles are considered conjugates. This is why angles are commonly used to represent parts of a whole, using the fan favorite, but often maligned, pie chart shown in Figure 3-6. The sum of the wedges makes a complete circle.

**Direction**

Direction is similar to angle, but instead of relying on two vectors joined at a point, direction relies on a single vector's orientation in a coordinate system. You can see which way is up, down, left, and right and everything in between.

This helps you determine slope, as shown in Figure 3-7. You can see increases, decreases, and fluctuations.

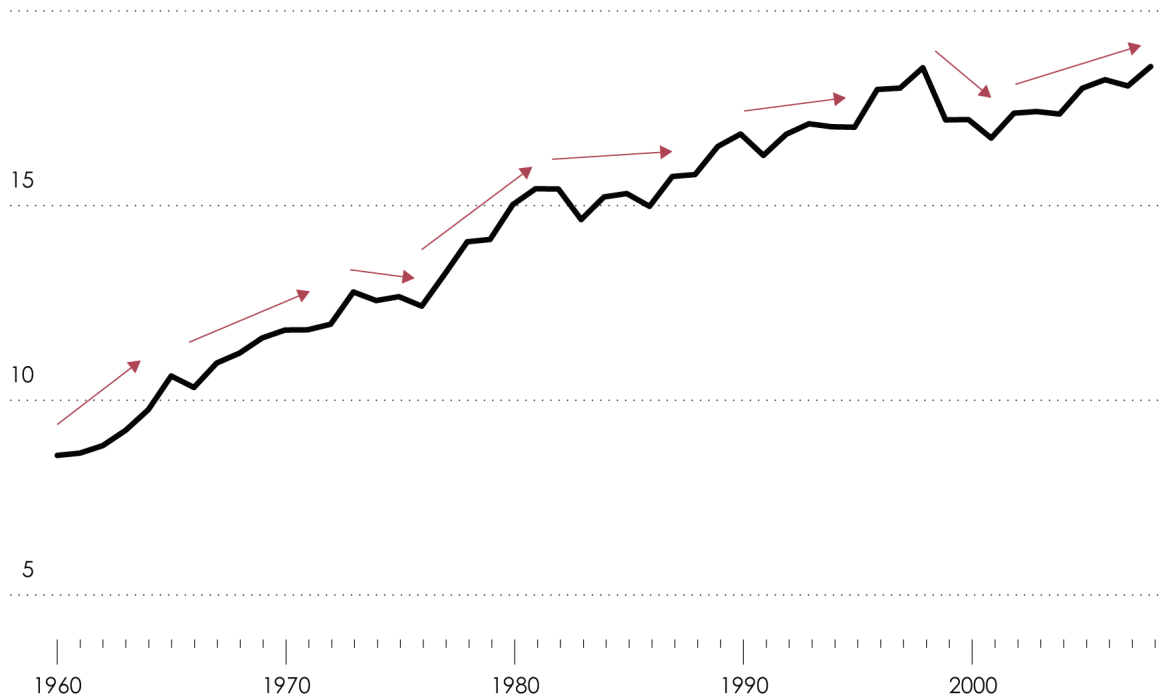
The amount of perceived change depends a lot on the scale, as shown in Figure 3-8. For example, you can make a small change in percentage look like

a lot by stretching out the scale. Likewise, you can make a big change look like a little by compressing the scale.

A rule of thumb is to scale your visualization so that direction fluctuates mostly around 45 degrees, but this is hardly a concrete rule. The best thing to do is to start with this suggestion and then adjust accordingly based on context. If a small change is significant, then it might be appropriate to stretch the scale so that you can see the shift. In contrast, if a small change is not significant, don't stretch out the scale just to make a shift look dramatic.

## Direction in a time series

20 metric tons of CO<sub>2</sub> per capita in Australia

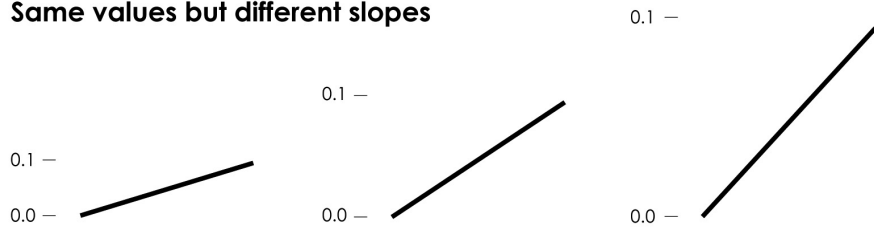


Source: The World Bank

**FIGURE 3-7** *Slope and time series*



### Same values but different slopes



**FIGURE 3-8** Same amount of change shown on varied scales

### Shapes

Shapes and symbols are commonly used with maps to differentiate categories and objects. Location on a map can be directly translated to the real world, so it makes sense to use icons to represent things in the real world. You might represent forests with trees or residential areas with houses.

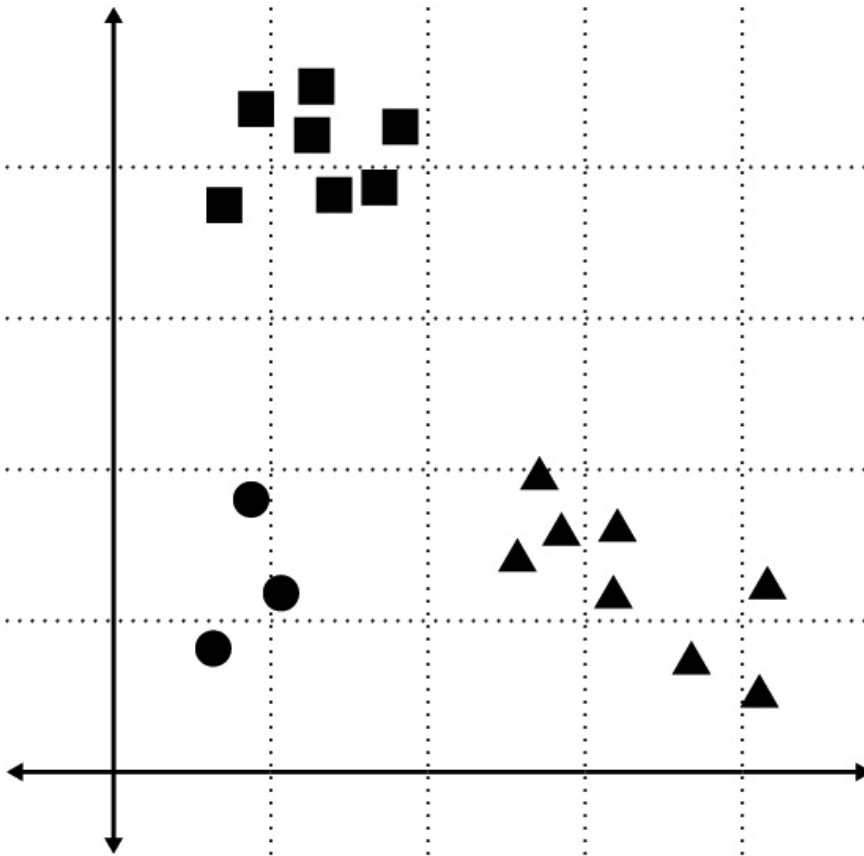
In a chart context, shapes to show variation are used less frequently than they used to be. When graphs were drawn with paper and a pencil and computers still worked with punch cards, symbols were an easy way to differentiate categories. For example, as shown in Figure 3-9, triangles and squares could be used in a scatterplot, which is quicker to draw than to switch between colored pencils and pens or fill a single shape with a solid or cross-hatched pattern.

Nevertheless, varied shapes can provide context that points alone can't, and it's typically not more difficult to try with your favorite software.

### Area and Volume

Bigger objects represent greater values. Like length, area and volume can be used to represent data with size, but with two and three dimensions, respectively. For the former, circles and rectangles are commonly used, and with the latter, cubes and sometimes spheres. You can also size more detailed icons and illustrations.

Be sure to mind how many dimensions you use. The most common mistake is to size a two- or three-dimensional object by only one dimension, such as height, but to maintain the proportions of all dimensions. This results in shapes that are too big and too small, which makes it impossible to fairly compare values.



**FIGURE 3-9** *Different shapes in scatterplot*

Say you use squares, shapes with two dimensions—width and height—to represent your data. The greater a value, the greater the area of a square, so if one value is 50 percent greater than another, you want the area of the square to be 50 percent greater than the other. However, if you increase the sides of the square by 50 percent instead of the area, which is what some software does by default, the larger square is too big. Instead of an increase in 50 percent, it's an increase of 125 percent. See the jump in difference in Figure 3-10.

You run into the same problem with three-dimensional objects, but the mistake is more pronounced. Increase the width, height, and depth of a cube by 50 percent, and the volume of the cube increases by approximately 238 percent.

Sizing by area

This is one unit.

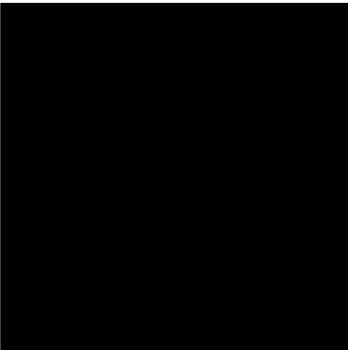


Four units sized by area



4 times the area as unit square

Four units *incorrectly* sized by side length



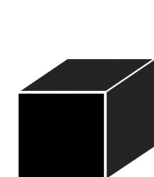
16 times the area as unit square

Sizing by volume

This is one unit.

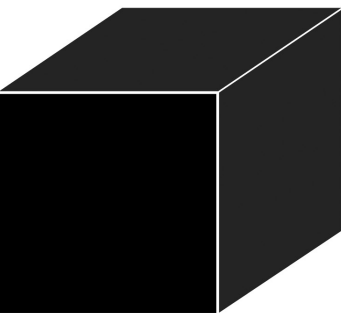


Four units sized by volume



4 times the volume as unit cube

Four units *incorrectly* sized by edge length



32 times the volume as unit square

FIGURE 3-10 Squares and cubes sized by different dimensions



## Color

Color as a visual cue can be spilt into two categories: hue and saturation. They can be used individually or in combination.

Color hue is what you usually just refer to as color. That's red, green, blue, and so on. Differing colors used together usually indicates categorical data, where each color represents a group. Saturation is the amount of hue in a color, so if your selected color is red, high saturation would be very red, and as you decrease saturation, it looks more faded. Used together, you can have multiple hues that represent categories, but each category can have varying scales.

Careful color selection can lend context to your data, and because there is no dependency on size or position, you can encode a lot of data at once. However, keep color blindness in mind if you want to make sure that everyone can interpret your graphics. Approximately 8 percent of men and 0.5 percent of women are red-green deficient, so when you encode your data only with those colors, this segment of your audience will have trouble decoding your visualization, if they can at all. Figure 3-11 shows how some shades are perceived by those who are color-deficient.

Does this mean you aren't allowed to use red and green in your graphics? No. You can combine visual cues so that everyone can make out differences, which you'll get to in a few sections.

**Note:** Something to think about: If someone is red-green deficient, how do they follow traffic signals that use red for stop and green for go? They note the order of the lights.



**FIGURE 3-11** Colors as perceived by those who have color vision deficiencies

Perception of Visual Cues

In 1985, William Cleveland and Robert McGill, then statistical scientists at AT&T Bell Laboratories, published a paper on graphical perception and methods. The focus of the study was to determine how accurately people read the visual cues above (excluding shapes), which resulted in a ranked list from most accurate to least accurate, as shown in Figure 3-12.

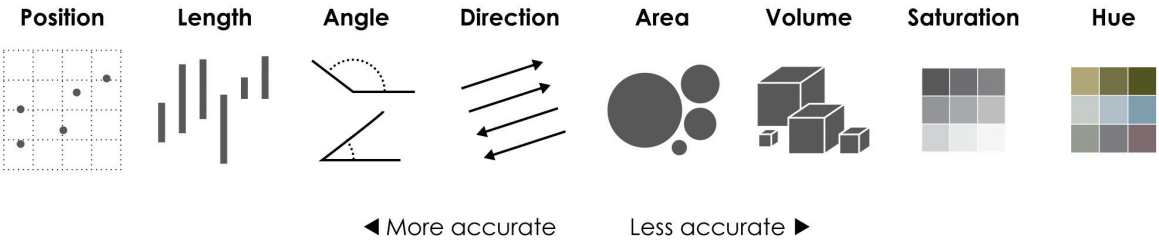


FIGURE 3-12 Visual cues ranked by Cleveland and McGill

A lot of visualization suggestions (and current research) stem from this list, which places bar charts above pie charts, heat maps at the bottom, and so on. This is sound advice, and you see more on this in Chapter 5, “Visualizing with Clarity,” but remember that this list doesn’t mean that dot plots are always better than bubble plots or that pie charts are evil.

Following this list blindly is an oversimplification of what visualization is. As you saw in the previous chapter, efficiency and exactness are not always the goal. That said, regardless of what you want to visualize data for, it’s good to know how well people can read your visual cues and what information they can extract. In other words, use these rankings as a guide rather than a rule book.

COORDINATE SYSTEMS

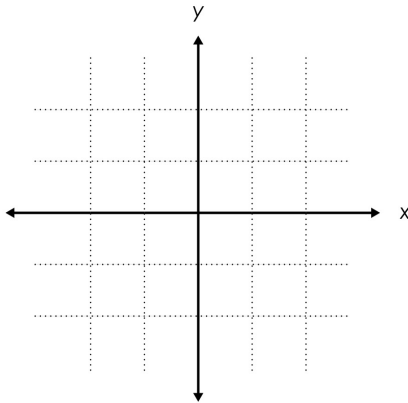
When you encode data, you eventually must place the objects somewhere. There’s a structured space and rules that dictate where the shapes and colors go. This is the coordinate system, which gives meaning to an x-y coordinate or a latitude and longitude pair. There are several systems, but as shown in Figure 3-13, there are three that cover most of your bases: Cartesian, polar, and geographic.

## Coordinate systems

There are a variety of them, from cylindrical to spherical, but these three will cover most of your bases.

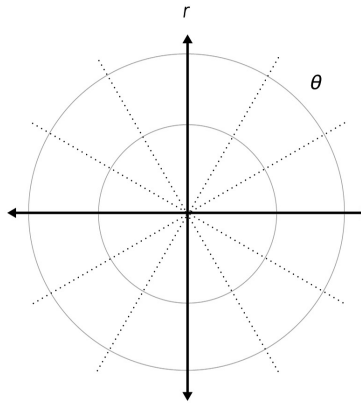
### Cartesian

If you've ever made a graph, the x- and y-coordinate system will look familiar to you.



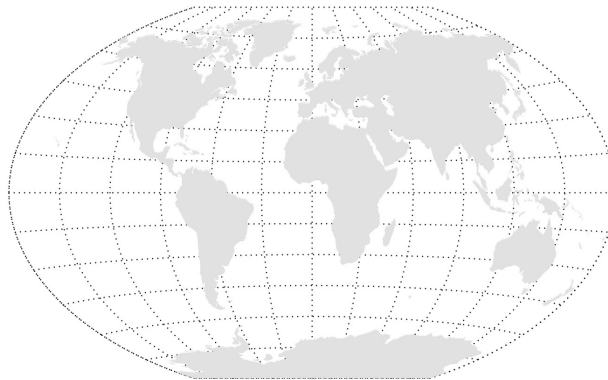
### Polar

Pie charts use this system. Coordinates are placed based on radius  $r$  and angle  $\theta$ .



### Geographic

Latitude and longitude are used to identify locations in the world. Because the planet is round, there are multiple projections to display geographic data in two dimensions. This one is the Winkel tripel.



**FIGURE 3-13** Commonly used coordinate systems



## Cartesian

The Cartesian coordinate system is the most commonly used one with charts. If you've made a traditional graph, such as a bar chart or a dot plot, you've used Cartesian coordinates.

You typically think of coordinates in the system as an  $x$  and  $y$  pair that is denoted as  $(x, y)$ . Two lines that are perpendicular to each other, and range from negative to positive, form the axes. The place the lines intersect is the origin, and the coordinate values indicate the distance from that origin. For example, the  $(x, y)$  pair at  $(0, 0)$  is at the intersection of the lines, and the  $(1, 2)$  pair is one unit away from the origin on the horizontal and two units away on the vertical.

To make this high school geometry flashback complete, you can find the distance between any two points, denoted as  $(x_1, y_1)$  and  $(x_2, y_2)$ , with the distance formula.

$$distance = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

You can also extend the Cartesian space to more than two dimensions. For example, a three-dimensional space would use a  $(x, y, z)$  triplet instead of just a  $(x, y)$  pair.

The takeaway is that you can describe geometric shapes using Cartesian coordinates, which makes it easier to draw in the space. From an implementation standpoint, the coordinate system enables you to encode values to paper or a computer screen.

## Polar

Made a pie chart? You've used the polar coordinate system, too. Although you might have used only the angle part and not the radius. Referring to Figure 3-13, the polar coordinate system consists of a circular grid, where the rightmost point is zero degrees. The greater the angle is, the more you rotate counterclockwise. The farther away from the circle you are, the greater the radius is.

Place yourself on the outer-most circle, and increase the angle. This rotates you counterclockwise toward the vertical line (or the  $y$ -axis if this were Cartesian coordinates), which is 90 degrees (that is, a right angle). Rotate one-quarter more, and you get to 180 degrees. Rotate back to where you started, and that's a 360-degree rotation. Your radius would be smaller if you rotated along the smaller circle.

This system is used less than the Cartesian coordinate system, but it can be useful in cases in which the angle or direction is important.

## Geographic

Location data has the added benefit of a connection to the physical world, which in turn lends instant context and a relationship to that point, relative to where you are. With a geographic coordinate system, you can map these points. Location data comes in many forms, but it's most commonly described as latitude and longitude, which are angles relative to the Equator and Prime Meridian, respectively. Sometimes elevation is also included.

Latitude lines run east and west, which indicates north and south position on a globe. Longitude lines run north and south and indicate the east and west position. Elevation can be thought of as a third dimension. Compared with Cartesian coordinates, latitude is like the horizontal axis, and longitude is like the vertical axis. That is, if you use a flat projection.

The tricky part about mapping the surface of Earth is that it's wrapped around a spherical mass, but you usually need to display it on a two-dimensional surface, like a computer screen. The variety of ways to do this are called projections, and as shown in Figure 3-14, each has its advantages and disadvantages.

When you project something that is three-dimensional onto a two-dimensional plane, some information is lost, whereas other information is preserved.

The Mercator projection, for example, preserves angles in local regions. It was created in the 16<sup>th</sup> century by cartographer Gerardus Mercator primarily for navigation on the seas and is still the most-used projection for online direction lookup. On the other hand, the Albers projection preserves area but distorts shape. So the projection you choose depends on what you want to focus on.

## SCALES

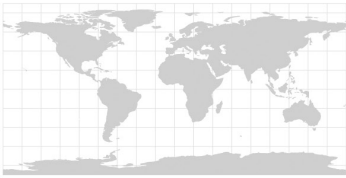
Whereas coordinate systems dictate the dimensions of a visualization, scale dictates where in those dimensions your data maps to. There's a variety of them, and you can even define your own scales based on mathematical functions, but most likely you'll rarely stray from the ones in Figure 3-15. These can be grouped into three categories: numeric, categorical, and time.

**FIGURE 3-14** (following page)  
*Map projections*

## Map projections

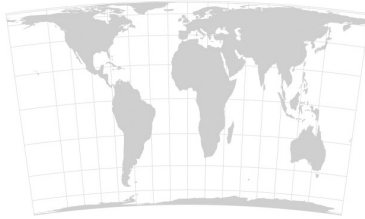
### Equirectangular

Typically used for thematic mapping, but doesn't preserve area or angle



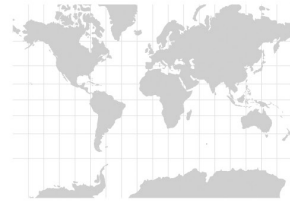
### Albers

Scale and shape not preserved; angle distortion is minimal



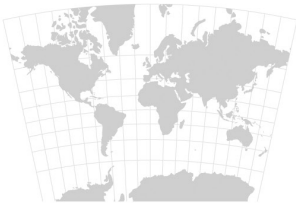
### Mercator

Preserves angles and shapes in small areas, making it good for directions



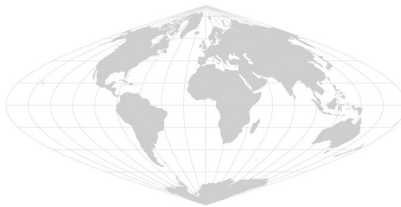
### Lambert conformal conic

Better for showing smaller areas and often used for aeronautical maps.



### Sinusoidal

Preserves area; useful for areas near the prime meridian



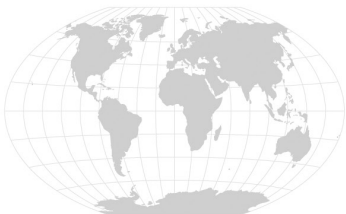
### Polyconic

Was often used to show US in the mid-1900s; little distortions in small areas near meridian



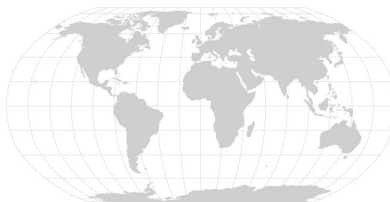
### Winkel Tripel

Minimized area, angle, and distance distortion; good choice for world map



### Robinson

A compromise between preserving areas and angles; good to show world map



### Orthographic

Representing a 3-D object in 2-D, need to rotate to area of interest

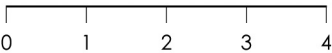


## Scales

Along with coordinate systems, they dictate where the shapes are placed and how objects are shaded.

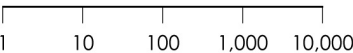
### Linear

Values are evenly spaced



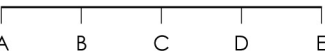
### Logarithmic

Focus on percent change



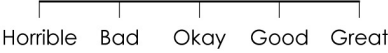
### Categorical

Discrete placement in bins



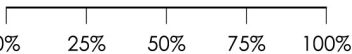
### Ordinal

Categories where order matters



### Percent

Representing parts of a whole



### Time

Units of months, days, or hours



FIGURE 3-15 Scales

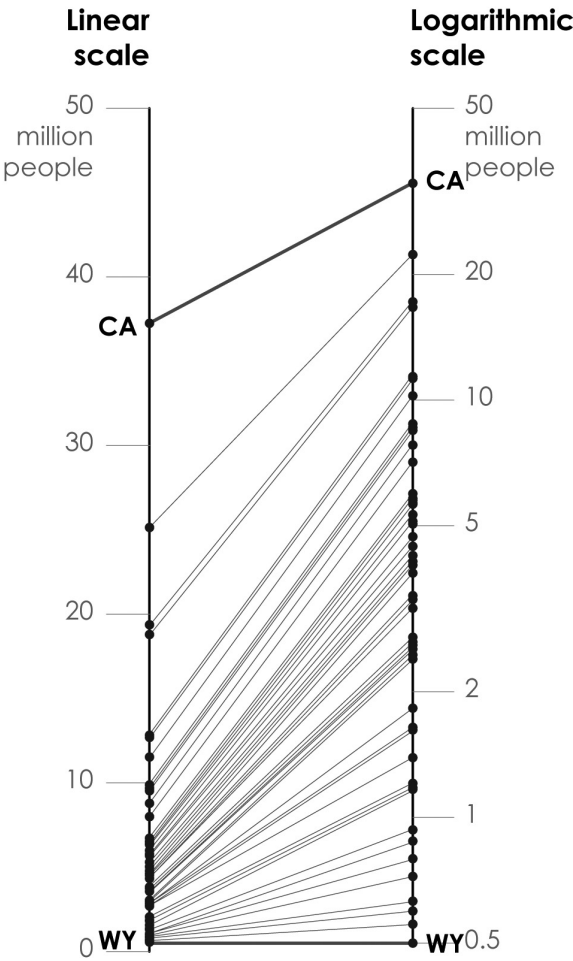
### Numeric

The visual spacing on a linear scale is the same regardless of where you are on the axis. So if you were to measure the distance between two points on the lower end of the scale, it'd be the same if they were at the high end of the scale.

On the other hand, a logarithmic scale condenses as you increase values. This scale is used less than the linear scale and is not as well understood or straightforward for those who don't regularly work with data, but it's useful if you're interested in percent differences more than you are raw counts or your data has a wide range.

For example, when you compare state populations in the United States, you deal with numbers from the hundreds of thousands up to the tens of millions. As of this writing, California has a population of approximately 38 million people, whereas Wyoming has a population of approximately 600,000. As shown

in Figure 3-16, with a linear scale, states with smaller populations are clustered on the bottom, and then a few states rest on top. It's easier to see points on the bottom with a logarithmic scale.



**FIGURE 3-16** *Linear versus logarithmic scale*

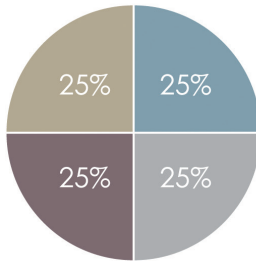
A percent scale is usually linear, but when it's used to represent parts of a whole, its maximum is 100 percent. As shown in Figure 3-17, the sum of all the parts is 100 percent. This seems obvious—that the sum of percentages in a pie chart, represented with wedges, should not exceed 100 percent—but the



mistake seems to come up occasionally. Sometimes it's due to mislabeling, but some people just aren't familiar with the concept.

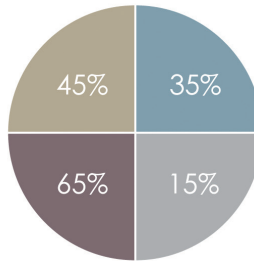
### Correct

The sum of the parts equals 100 percent.



### Wrong

The sum of the parts is more than 100 percent.



Also, the labels don't match the wedge sizes.

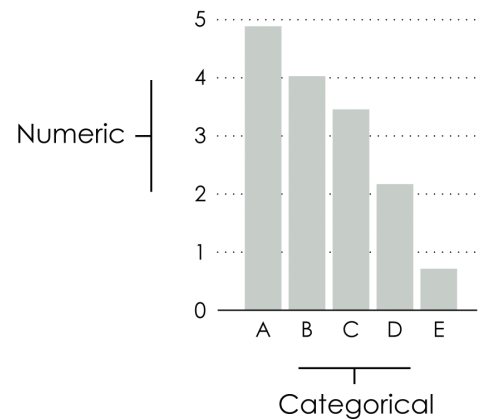
**FIGURE 3-17** *Incorrect and correct pie charts*

### Categorical

Data doesn't always need to be numeric. It can be categorical, such as people's cities of residence or the political parties of government officials. A categorical scale provides visual separation for these different groups and often works with a numeric scale. A bar plot for example, can use a categorical scale on the horizontal axis and a numeric scale on the vertical to show counts or measurements for different groups, as shown in Figure 3-18.

Spacing between each category is arbitrary because it does not depend on a numeric value, but it is typically adjusted to increase clarity, which is discussed in Chapter 6, "Visualizing with Clarity."

Ordering should be used in the context of the data. Although this can also be arbitrary, for an ordinal scale that uses categories, order of course matters. If your data is a categorical ranking on movies that ranges from horrible to great, then it makes sense to keep that order visually, which makes it easier to compare and judge quality.



**FIGURE 3-18** *Numeric and categorical scales on a bar plot*

Time

Time is a continuous variable, which lets you plot temporal data on a linear scale, but you can divide it into categories such as months or days of the week, which lets you visualize it as a discrete variable. Also, it cycles, as shown in Figure 3-19. There’s always another noon time, Saturday, and January.

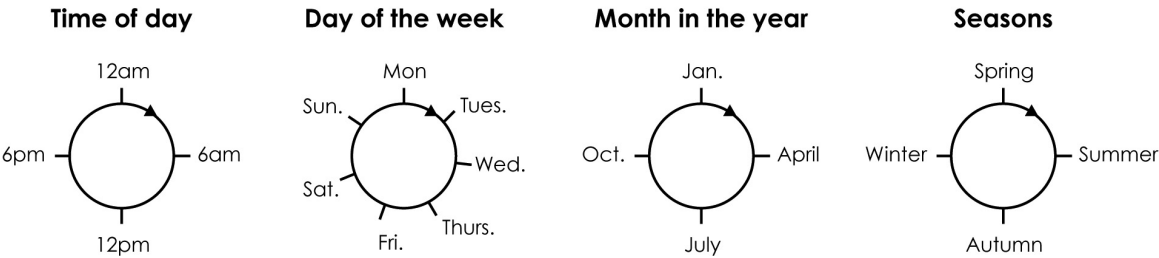


FIGURE 3-19 Time cycles

You saw this in Chapter 1, which showed fatal crashes over time, by year, by month, by day, and by hour. Data was plotted continuously in these cases. However, aggregates by time of day, day of the week, and month (over multiple years) showed a different picture.

When communicating data to an audience, the time scale, like geographic maps, gives you an advantage of lending a reader connection because time is a part of everyday life. You feel and experience time internally and through your clocks and calendars, and as the sun rises and sets.

CONTEXT

Context (information that lends to better understanding the who, what, when, where, and why of your data) can make the data clearer for readers and point them in the right direction. At the least, it can remind you what a graph is about when you come back to it a few months later.

Sometimes context is explicitly drawn, and other times it’s implied through the medium. For example, as shown in Figure 3-20, designers Matt Robinson

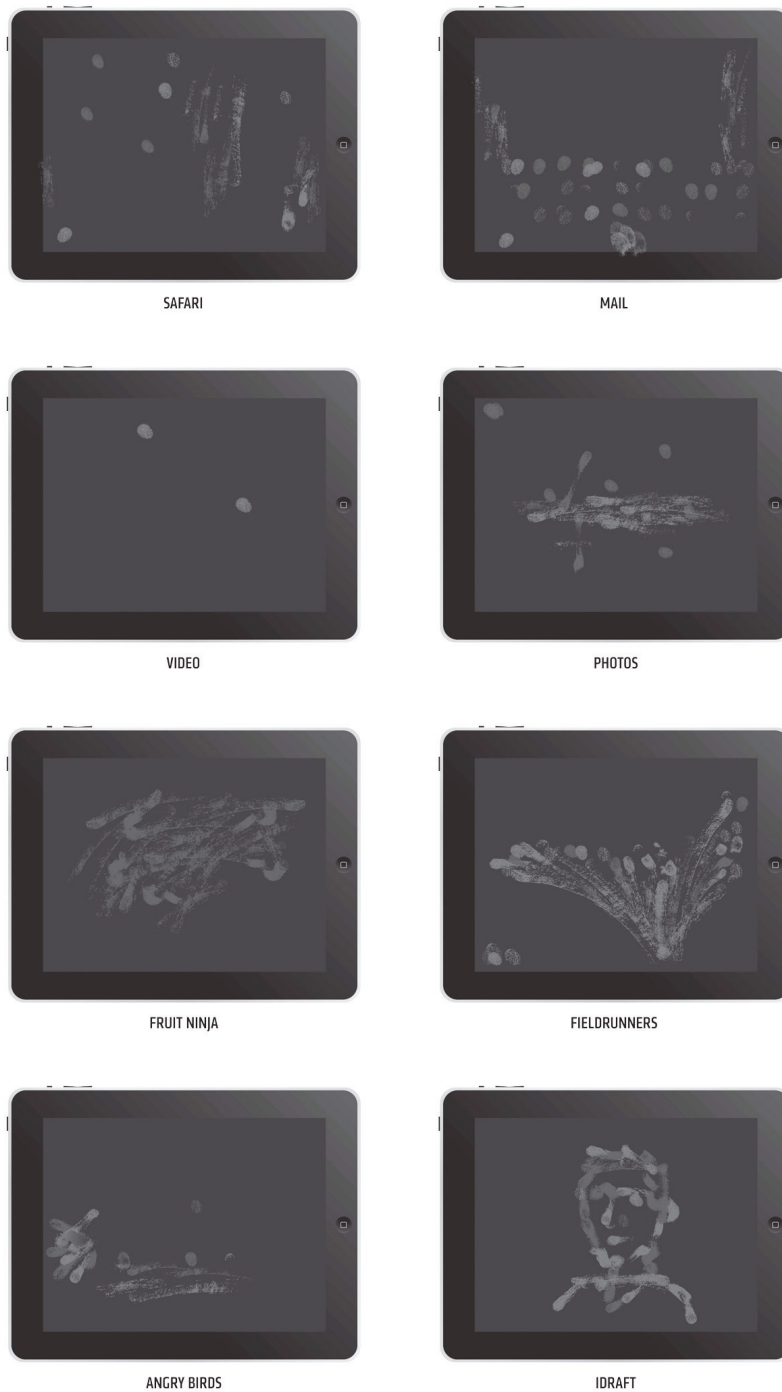
and Tom Wigglesworth drew “sample” on a wall, with ballpoint pens and different typefaces. Because ink usage varies by typeface, each pen had a different amount of ink left, which made for an interesting bar graph. There’s no need to label the numeric axis because it’s implied by the pens and their ink.

Designer George Kokkinidis approached iPad usage in a similar way; however, as shown in Figure 3-21, instead of comparing remaining ink, he looked at fingerprint traces while he used different apps. For example, in Mail, he typed messages most of the time, so the keyboard pattern is most evident, with some scrolling on the side. In contrast, most interaction is in the bottom-left corner for the game Angry Birds.

Of course, you can’t always draw on familiar physical objects for context, so you must provide familiarity and a sense of scale in other ways. The easiest and most straightforward way is to label your axes and specify units of measure, or provide a description that tells others what each visual cue represents. Otherwise, when the data is abstracted, there’s no way to decode the shapes, sizes, and colors, and you might as well show an amorphous blob.



**FIGURE 3-20** Measuring Type (2010) by Matt Robinson and Tom Wigglesworth, <http://datafl.ws/27m>



**FIGURE 3-21** Remnants of a Disappearing UI (2010) by George Kokkinidis, <http://datafl.ws/27n>

A descriptive title is a small but easy thing you can create to set up readers for what they're about to look at. Imagine you produce a time series plot for gas prices that shows an upward trend. You could just title it "Gas Prices" and that would be a fair title. That's what it is, but you could also title it "Rising Gas Prices," which says what data is used and what is shown. You could also include lead-in text underneath the title that describes fluctuations or by how much gas prices rose.

Your choice of visual cues, a coordinate system, and scale can implicitly provide context. Bright, cheery, and contrasting colors says something different than dark, neutral, and blending colors. Similarly, a geographic coordinate system places you within the context of physical space, whereas an x-y plot using Cartesian coordinates keeps you within a virtual space. A logarithmic scale could suggest a focus on percentage changes and reduce focus on absolute values.

This is why it's important to pay attention to software defaults.

Programs are designed to be flexible and fast and they work outside the context of the data. This is great to draw a visualization base and explore your data, but it's up to you to make the right decisions along the way and to make the computer output something for humans. This comes partly from knowing how you perceive geometry and colors, but mostly it comes from practice and the experience gained from seeing a lot of data and evaluating how others, who aren't familiar with your data, interpret your work. Common sense also goes a long way.

**Note:** The people who visualize data best got to that level because they examined and visualized a lot of data. They gained experience with each graph made. Reading books will inform better decisions, but it's not until you put what you learn into practice when you really improve.

## PUTTING IT TOGETHER

You know what ingredients are available. Now it's time to cook the meal. Viewed separately, the visualization components aren't that useful because they are just bits of geometry floating in an empty space without context. However, when you put the components together, you get a complete visualization worth looking at.

For example, what do you get when you use length as a visual cue, a Cartesian coordinate system, and a categorical scale on the horizontal axis and a linear scale on the vertical? You get a bar chart. Use position with a geographic coordinate system, and you get points on a map.

What do you get when you use a polar coordinate system with the area as the visual cue, a percentage scale on the radius, and a time scale on the rotation? That’s a polar area diagram. The most famous one, as shown in Figure 3-22, is Florence Nightingale’s chart that visualizes deaths from treatable diseases over time.

*On the Origin of Species: The Preservation of Favoured Races*, by designer and developer Ben Fry, uses color and length, Cartesian coordinates, and a linear scale, as shown in Figure 3-23. The interactive and animated visualization shows how Charles Darwin’s theory of evolution changed through six editions. The gray blocks represent the original text, and each subsequent color represents a revision in an edition, so you can see what changed and by how much.

DIAGRAM OF THE CAUSES OF MORTALITY  
IN THE ARMY IN THE EAST.

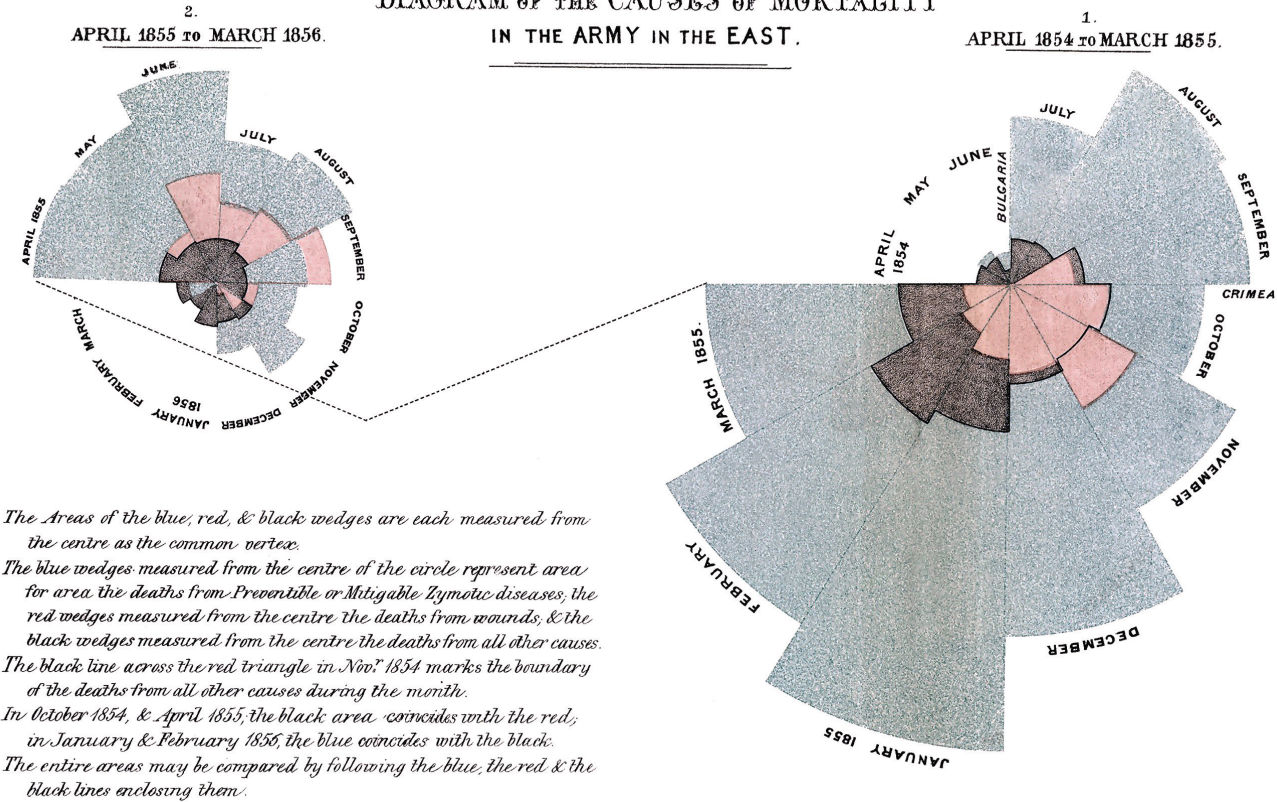


FIGURE 3-22 Diagram of the Causes of Mortality in the Army in the East (1858) by Florence Nightingale



# ON THE ORIGIN OF SPECIES *The Preservation of Favoured Races*

Reset Pause Slow Fast



First Edition (1859) Second Edition (1860) Third Edition (1861) Fourth Edition (1866) Fifth Edition (1869) Sixth Edition (1872)

FIGURE 3-23 On the Origin of Species: The Preservation of Favoured Races (2009) by Ben Fry, <http://benfry.com/traces/>



In the deaths chart shown in Figure 3-24, from the *Statistical Atlas of the United States* published in 1874, length is used to show the distribution of deaths for each state, by age and gender. The horizontal axis on each plot represents the number of deaths on a linear scale, and the vertical axis represents numeric categories that represent age groups.

Figure 3-25 shows generalized combinations, which cover common visualization types such as the line chart, bubble plot, and choropleth map. The key is learning how each component fits together—how each ingredient can complement and enhance others—to make something more useful than the separate parts.

Now try to fit components together, starting with the data and then building on that foundation. Figure 3-26 is a data table from the United States Census Bureau that shows educational attainment (high school graduate or more, bachelor’s degree or more, and advanced degree or more) by state, in 1990, 2000, and 2009. The values are percentages for people 25 years old and over. These are important bits about the data that you need to know before actually looking at the data.

The “or more” for each column means you can’t just add the values from each column because there’s overlap between them. If you want to make a pie chart that shows the values of each column, you must do some math. For example, the United States estimate for people with a high school degree (or equivalent) or more is 75.2 percent. Subtract those with a bachelor’s degree or more, 20.3 percent, to get rid of the “or more” part of the high school value, which gives you 54.9 percent of people with only a high school degree.

It’s also useful to know the sample population. If it were everyone in America, the percentages would be lower, or if for some odd reason the sample was those under 18, the percentages for an advanced degree or more would represent a tiny group of people who skipped or advanced quickly through elementary and high school.

So you have the most important part of any visualization: the data. There are nine columns, spread out over 3 years and three subcategories, plus one more column for state names, so you can visualize the data on multiple dimensions. You might want to focus on educational attainment in 2009, in which case, a few bar charts, as shown in Figure 3-27, could work.

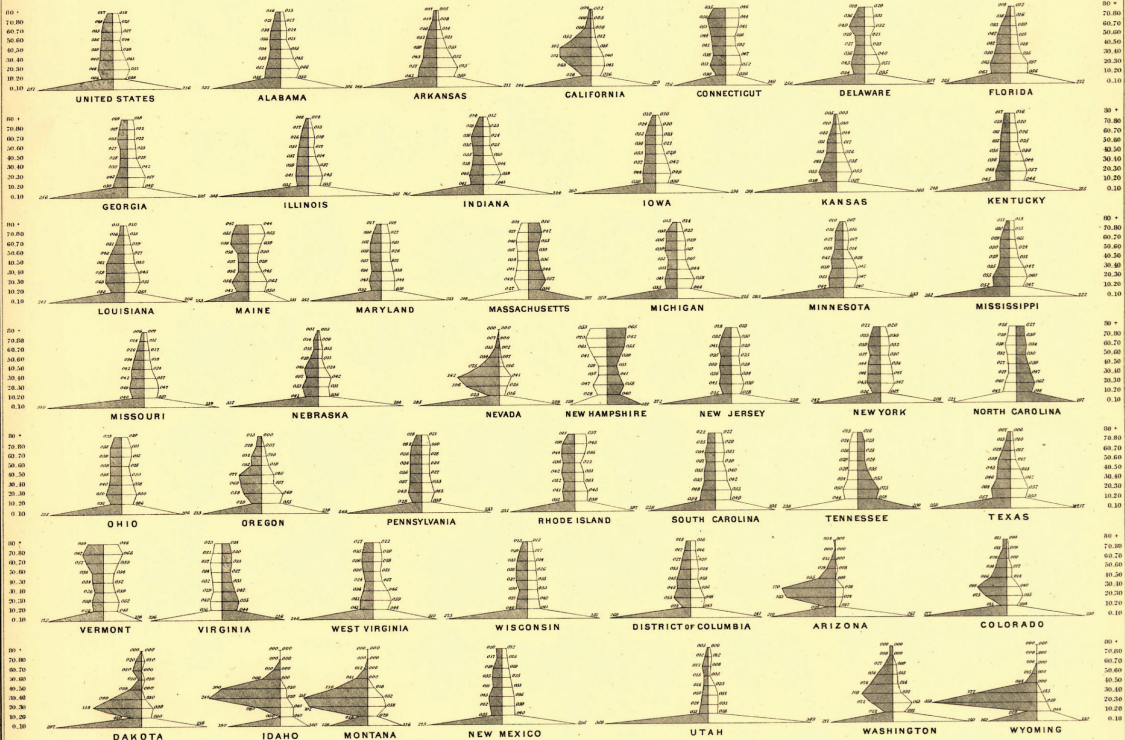
**FIGURE 3-24** (facing page)  
Chart showing the distributions  
of deaths, based on United States  
census of 1870 by Francis A. Walker

CHART  
SHOWING THE DISTRIBUTION BY AGE AND SEX  
OF 1000  
**DEATHS**  
OCCURRING DURING THE CENSUS YEAR ENDING JUNE 1 187,  
compiled from the Returns of Mortality at the Sixth Census 1870  
FRANCIS A. WALKER.

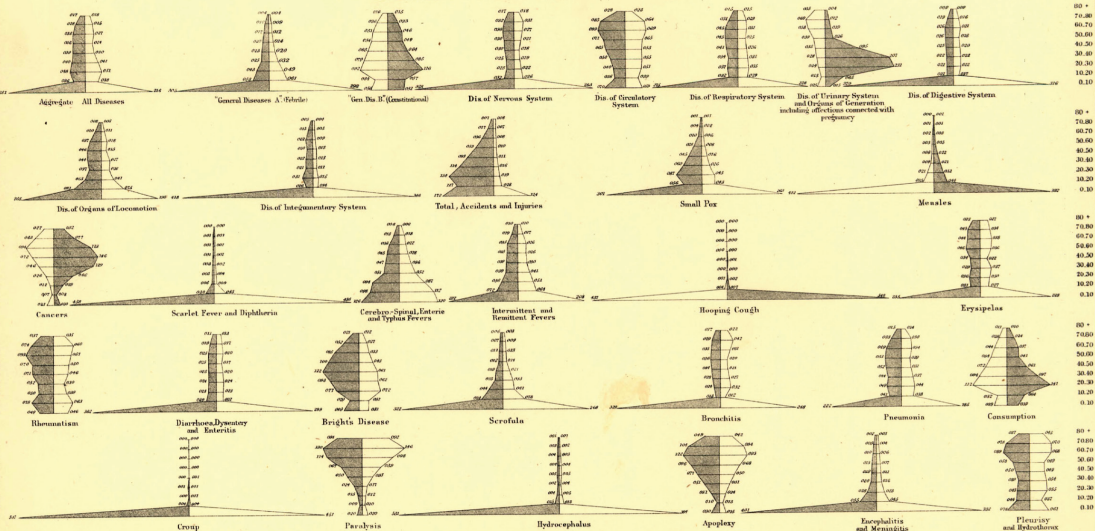
The total number of Deaths in each State or from each special Disease or Group of Diseases as reported in the Census is reduced to thousandths, and the number of thousandths of each sex in each decade or left, is represented by the distance measured on the horizontal line, sevenfold from the perpendicular base line.

The males are on the left of the perpendicular base line and the females on the right.  
The lowest horizontal line represents the deaths in the first decade, under ten years of age, and the highest the deaths over eighty years.  
The sex which preponderates is shaded.

1. FOR THE UNITED STATES AND FOR THE SEVERAL STATES AND TERRITORIES.



2. FOR GROUPS OF DISEASES AND CERTAIN SPECIAL DISEASES.



# Visual cues

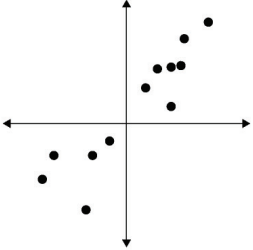

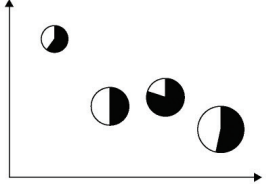
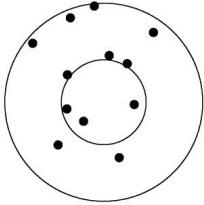
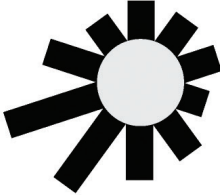
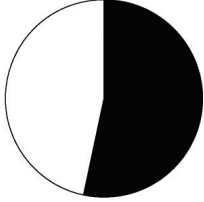



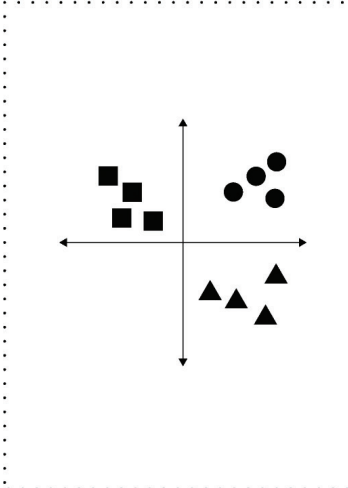
	Position	Length	Angle
Coordinate systems			
Cartesian			
Polar			
Geographic			

FIGURE 3-25 Visualization component combinations

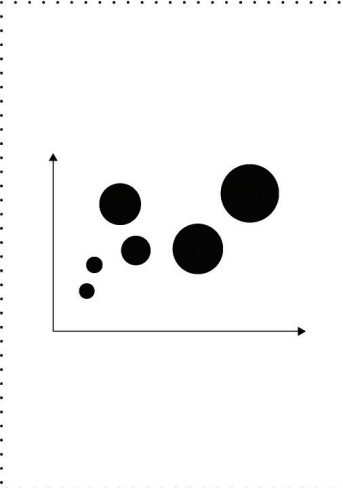
Direction



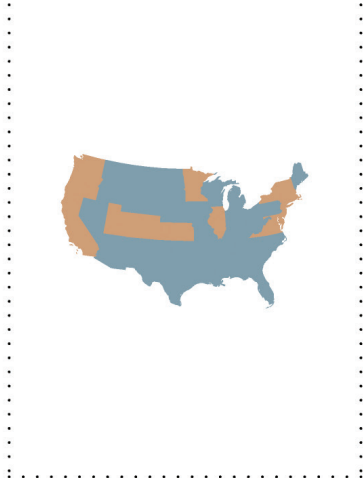
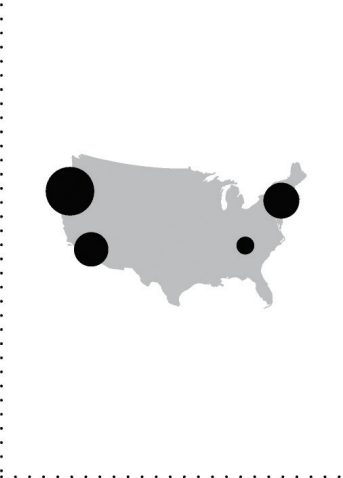
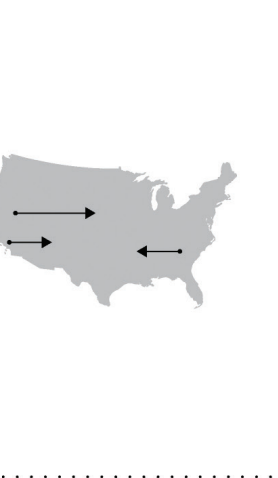
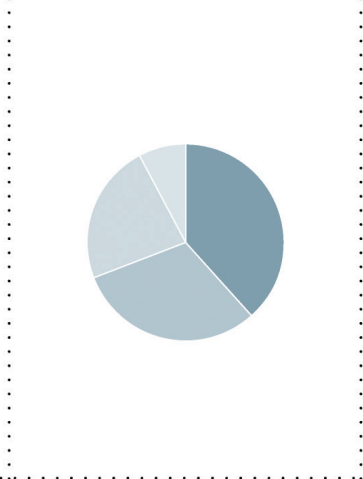
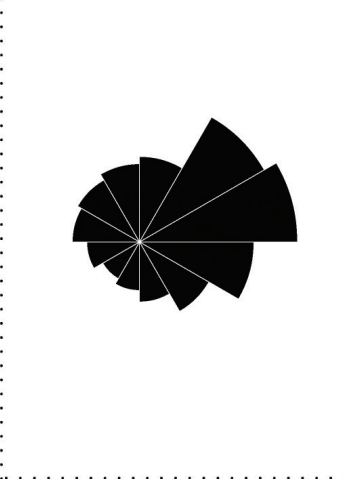
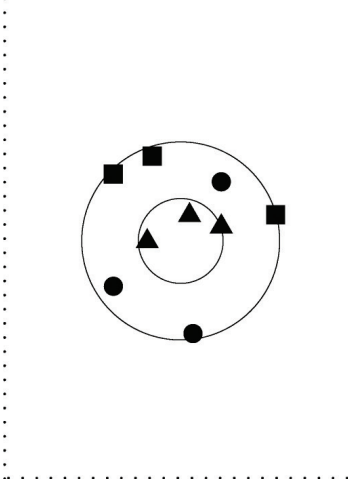
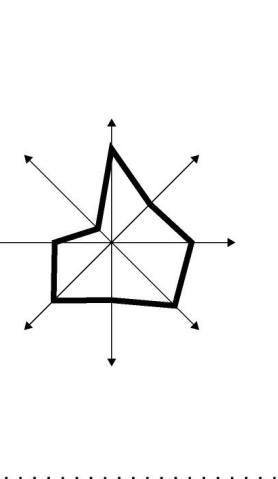
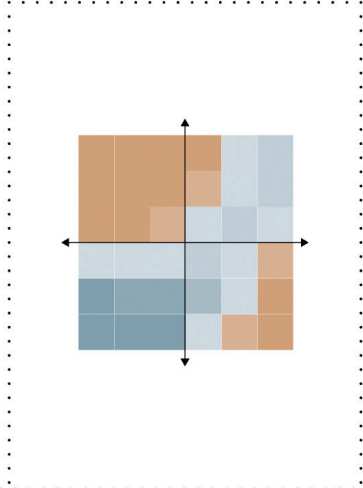
Shapes



Area or Volume



Color





**Table 233. Educational Attainment by State: 1990 to 2009**

[In percent. 1990 and 2000 as of April. 2009 represents annual averages for calendar year. For persons 25 years old and over. Based on the 1990 and 2000 Census of Population and the 2009 American Community Survey, which includes the household population and the population living in institutions, college dormitories, and other group quarters. See text, Section 1 and Appendix III. For margin of error data, see source]

State	1990			2000			2009		
	High school graduate or more	Bachelor's degree or more	Advanced degree or more	High school graduate or more	Bachelor's degree or more	Advanced degree or more	High school graduate or more	Bachelor's degree or more	Advanced degree or more
<b>United States . . . . .</b>	<b>75.2</b>	<b>20.3</b>	<b>7.2</b>	<b>80.4</b>	<b>24.4</b>	<b>8.9</b>	<b>85.3</b>	<b>27.9</b>	<b>10.3</b>
Alabama . . . . .	66.9	15.7	5.5	75.3	19.0	6.9	82.1	22.0	7.7
Alaska . . . . .	86.6	23.0	8.0	88.3	24.7	8.6	91.4	26.6	9.0
Arizona . . . . .	78.7	20.3	7.0	81.0	23.5	8.4	84.2	25.6	9.3
Arkansas . . . . .	66.3	13.3	4.5	75.3	16.7	5.7	82.4	18.9	6.1
California . . . . .	76.2	23.4	8.1	76.8	26.6	9.5	80.6	29.9	10.7
Colorado . . . . .	84.4	27.0	9.0	86.9	32.7	11.1	89.3	35.9	12.7
Connecticut . . . . .	79.2	27.2	11.0	84.0	31.4	13.3	88.6	35.6	15.5
Delaware . . . . .	77.5	21.4	7.7	82.6	25.0	9.4	87.4	28.7	11.4
District of Columbia . . . . .	73.1	33.3	17.2	77.8	39.1	21.0	87.1	48.5	28.0
Florida . . . . .	74.4	18.3	6.3	79.9	22.3	8.1	85.3	25.3	9.0
Georgia . . . . .	70.9	19.3	6.4	78.6	24.3	8.3	83.9	27.5	9.9
Hawaii . . . . .	80.1	22.9	7.1	84.6	26.2	8.4	90.4	29.6	9.9
Idaho . . . . .	79.7	17.7	5.3	84.7	21.7	6.8	88.4	23.9	7.5
Illinois . . . . .	76.2	21.0	7.5	81.4	26.1	9.5	86.4	30.6	11.7
Indiana . . . . .	75.6	15.6	6.4	82.1	19.4	7.2	86.6	22.5	8.1
Iowa . . . . .	80.1	16.9	5.2	86.1	21.2	6.5	90.5	25.1	7.4
Kansas . . . . .	81.3	21.1	7.0	86.0	25.8	8.7	89.7	29.5	10.2
Kentucky . . . . .	64.6	13.6	5.5	74.1	17.1	6.9	81.7	21.0	8.5
Louisiana . . . . .	68.3	16.1	5.6	74.8	18.7	6.5	82.2	21.4	6.9
Maine . . . . .	78.8	18.8	6.1	85.4	22.9	7.9	90.2	26.9	9.6
Maryland . . . . .	78.4	26.5	10.9	83.8	31.4	13.4	88.2	35.7	16.0
Massachusetts . . . . .	80.0	27.2	10.6	84.8	33.2	13.7	89.0	38.2	16.4
Michigan . . . . .	76.8	17.4	6.4	83.4	21.8	8.1	87.9	24.6	9.4
Minnesota . . . . .	82.4	21.8	6.3	87.9	27.4	8.3	91.5	31.5	10.3
Mississippi . . . . .	64.3	14.7	5.1	72.9	16.9	5.8	80.4	19.6	7.1
Missouri . . . . .	73.9	17.8	6.1	81.3	21.6	7.6	86.8	25.2	9.5
Montana . . . . .	81.0	19.8	5.7	87.2	24.4	7.2	90.8	27.4	8.3
Nebraska . . . . .	81.8	18.9	5.9	86.6	23.7	7.3	89.8	27.4	8.8
Nevada . . . . .	78.8	15.3	5.2	80.7	18.2	6.1	83.9	21.8	7.6
New Hampshire . . . . .	82.2	24.4	7.9	87.4	28.7	10.0	91.3	32.0	11.2
New Jersey . . . . .	76.7	24.9	8.8	82.1	29.8	11.0	87.4	34.5	12.9
New Mexico . . . . .	75.1	20.4	8.3	78.9	23.5	9.8	82.8	25.3	10.4
New York . . . . .	74.8	23.1	9.9	79.1	27.4	11.8	84.7	32.4	14.0
North Carolina . . . . .	70.0	17.4	5.4	78.1	22.5	7.2	84.3	26.5	8.8
North Dakota . . . . .	76.7	18.1	4.5	83.9	22.0	5.5	90.1	25.8	6.7
Ohio . . . . .	75.7	17.0	5.9	83.0	21.1	7.4	87.6	24.1	8.8
Oklahoma . . . . .	74.6	17.8	6.0	80.6	20.3	6.8	85.6	22.7	7.4
Oregon . . . . .	81.5	20.6	7.0	85.1	25.1	8.7	89.1	29.2	10.4
Pennsylvania . . . . .	74.7	17.9	6.6	81.9	22.4	8.4	87.9	26.4	10.2
Rhode Island . . . . .	72.0	21.3	7.8	78.0	25.6	9.7	84.7	30.5	11.7
South Carolina . . . . .	68.3	16.6	5.4	76.3	20.4	6.9	83.6	24.3	8.4
South Dakota . . . . .	77.1	17.2	4.9	84.6	21.5	6.0	89.9	25.1	7.3
Tennessee . . . . .	67.1	16.0	5.4	75.9	19.6	6.8	83.1	23.0	7.9
Texas . . . . .	72.1	20.3	6.5	75.7	23.2	7.6	79.9	25.5	8.5
Utah . . . . .	85.1	22.3	6.8	87.7	26.1	8.3	90.4	28.5	9.1
Vermont . . . . .	80.8	24.3	8.9	86.4	29.4	11.1	91.0	33.1	13.3
Virginia . . . . .	75.2	24.5	9.1	81.5	29.5	11.6	86.6	34.0	14.1
Washington . . . . .	83.8	22.9	7.0	87.1	27.7	9.3	89.7	31.0	11.1
West Virginia . . . . .	66.0	12.3	4.8	75.2	14.8	5.9	82.8	17.3	6.7
Wisconsin . . . . .	78.6	17.7	5.6	85.1	22.4	7.2	89.8	25.7	8.4
Wyoming . . . . .	83.0	18.8	5.7	87.9	21.9	7.0	91.8	23.8	7.9

Source: U.S. Census Bureau, 1990 Census of Population, CPH-L-96; 2000 Census of Population, P37. "Sex by Educational Attainment for the Population 25 Years and Over"; 2009 American Community Survey, R1501, "Percent of Persons 25 Years and Over Who Have Completed High School (Includes Equivalency)," R1502, "Percent of Persons 25 Years and Over Who Have Completed a Bachelor's Degree," and R1503, "Percent of Persons 25 Years and Over Who Have Completed an Advanced Degree," <<http://factfinder.census.gov/>>, accessed February 2011.

# Educational attainment in 2009

At least...

High school graduate

Bachelor's degree

Advanced degree



FIGURE 3-27 Bar charts on educational attainment

This is practically a direct translation of the last three columns in the table. Each row represents the values for a state, and each column is a level of attainment. Each bar chart has its own linear scale, but the increments are spaced equally and start at zero percent. States are sorted by estimated percent of people with a high school diploma or equivalent, in descending order, rather than alphabetically, like in the table. Instead of giving the national average its own row, it's presented as a vertical dotted line to provide a sense of low and high. Color hue—gray, light blue, and blue—is used to indicate three separate estimates.

Break it down. That's length (bars), color (each bar chart), and position (lines for national averages) as visual cues, a Cartesian coordinate system, linear scales for each of the bar charts, and a categorical scale for the sorted states. The title and subtitles provide context for what the data is about.

If you are more interested in the changes between 2000 and 2009 than you are just the 2009 percentages, Figure 3-28 shows a few options that shift focus. Length and position are still used, as well as a linear scale on the horizontal axis and a categorical scale on the vertical. However, the context and layout are different than the bar charts. Some other visual cues are also incorporated.

An open circle represents the high school attainment in 2000 for each state, and the solid circles represent the same for 2009. The dots are placed in the same position vertically, and a line is used to connect the two dots. The longer the line is, the greater the change, by percentage points, was from 2000 to 2009.

The shift from open circle to closed circle provides a sense of direction. In this example, high school attainment in all states improved, so your eyes always shift from left to right, but if attainment decreased in one of the states, you could use the same visual cue. For example, if there were a decrease from 80 percent to 70 percent, the solid dot would be on the left of the open one. You can also use arrows if you want to highlight direction more prominently. All states showed increases in this example, though, so a focus on the magnitude of the changes and the values of the endpoints was more appropriate.

You can see how a change in sorting can shift focus. States are sorted alphabetically in the first chart, and the lack of visual order makes it more challenging to make comparisons. You can see the increases and it's easy to find a state of interest, but as an overall picture, you don't get much.



In contrast, the second chart shows the same data ordered instead by the highest percentage of attainment in 2009. It starts with Wyoming and goes down to Texas. This focuses on the more recent estimates, whereas still making it easy to pick out the values for 2000 because generally speaking, states with higher percentages in 2009 were higher in the rankings in 2000, too. That said, you can also sort by the 2000 estimates and move the labels to the left to shift focus in this direction.

Finally, the chart on the far right introduces color as a visual cue. This is the same as the second chart that sorts by 2009 estimates, but color is used to highlight states that increase the most by percentage. The District of Columbia, which albeit isn't a state, had the greatest percentage increase, so it is shown in black. The lower the increase, the lighter the states are shown. States in between are shown with varying shades of green. So if you look at the individual components of this chart, you get length, position, direction, and color used as visual cues; it uses a Cartesian coordinate system; and a linear numeric scale is used on the horizontal, with a categorical scale on the vertical.

You don't have to stop here. As shown in Figure 3-29, position and direction can be used differently to show the increases from 2000 to 2009. Unlike the previous charts, states are plotted on a linear scale that represents high school attainment instead of on a categorical scale. Values are categorized by year on the horizontal. This is essentially a couple of ticks on a time series plot. If you were to show years in between, there would be more than two categories on the horizontal axis. In any case, like in a time series plot, a greater slope from point to point means a greater rate of change.

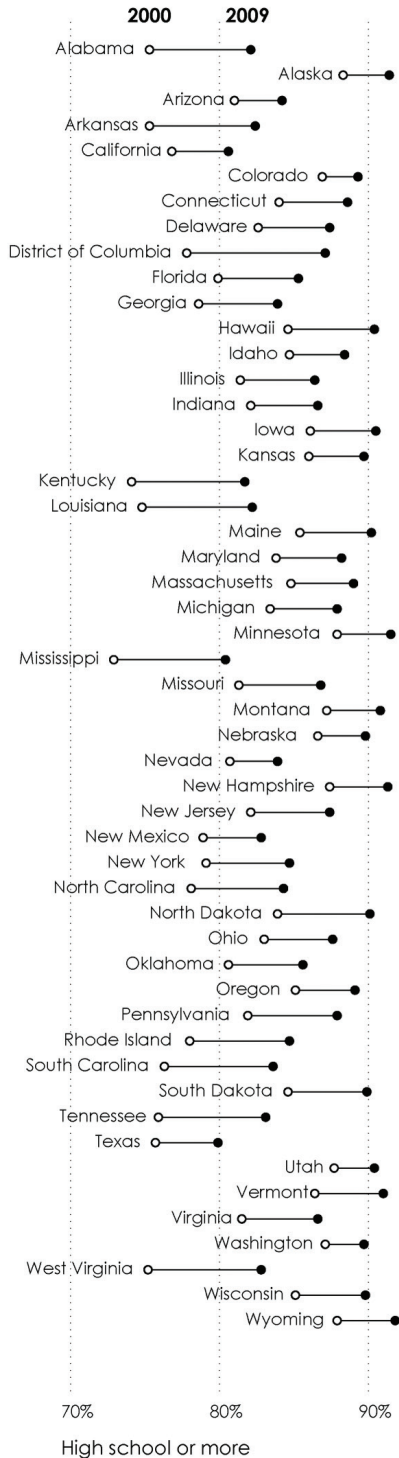
The chart on the right uses the same geometry as the one on the left, and uses color to represent regions in the United States. So although you see improvement with all states, you also see a lot of the states in the South toward the bottom of the scale and Midwest and West states more toward the top. Although, as is usually the case with real data, there are exceptions, such as California in the West that is toward the bottom and Maryland that is in the South is higher up.

Generally speaking though, the higher the attainment in 2000, the higher the attainment was in 2009. This is obvious in Figure 3-30, which uses position as a visual cue and linear scales on both axes.

**FIGURE 3-28** (following page)  
*Change in high school educational attainment between 2000 and 2009*

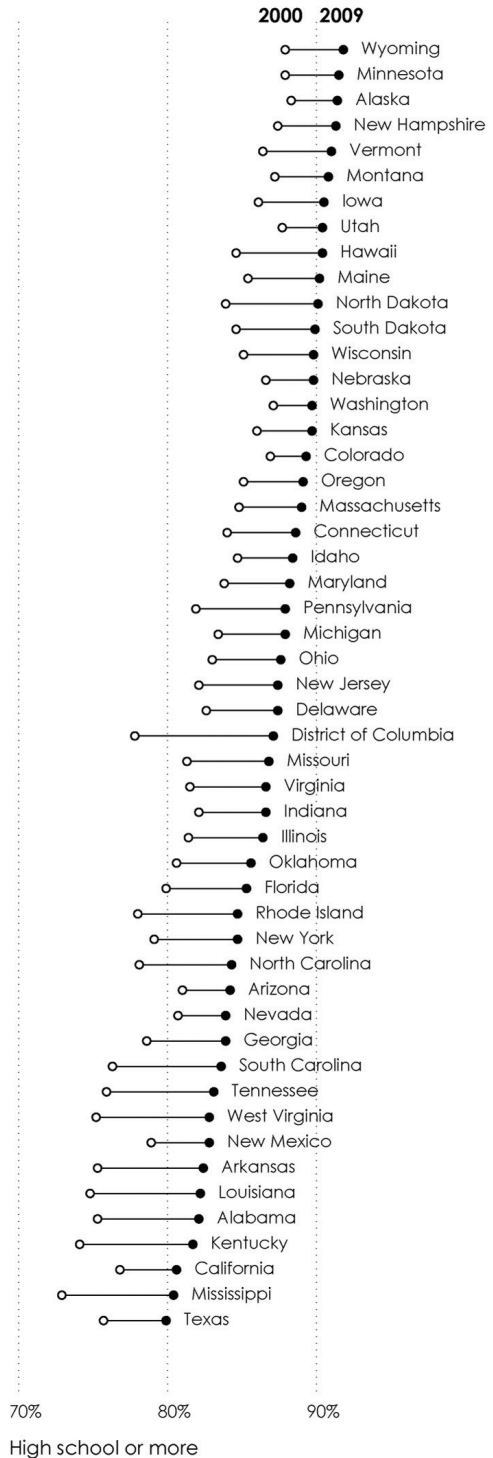
## Alphabetical

Values look scattered



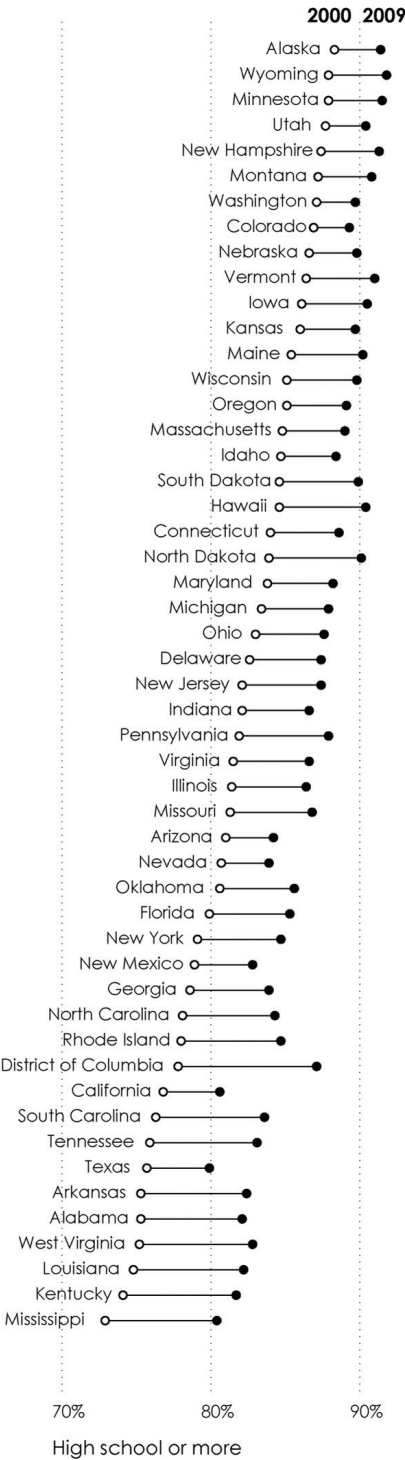
## Greatest to least, endpoint

Focus on most recent numbers



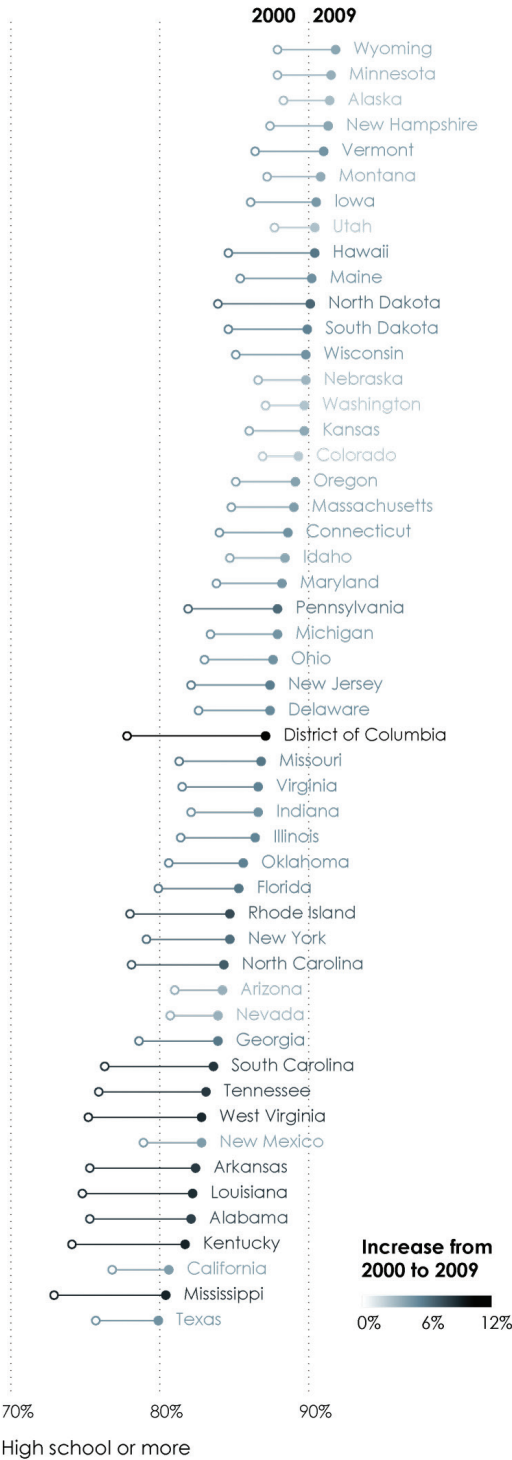
Greatest to least, startpoint

Focus on past data

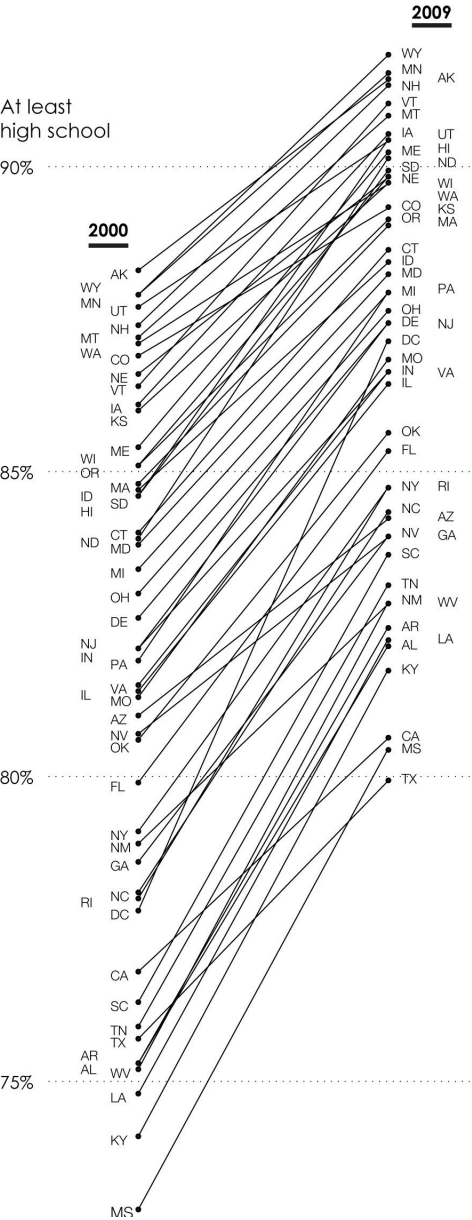


Greatest to least, endpoint + Color

Reference to percent increase



Position + Direction



Position + Direction + Color

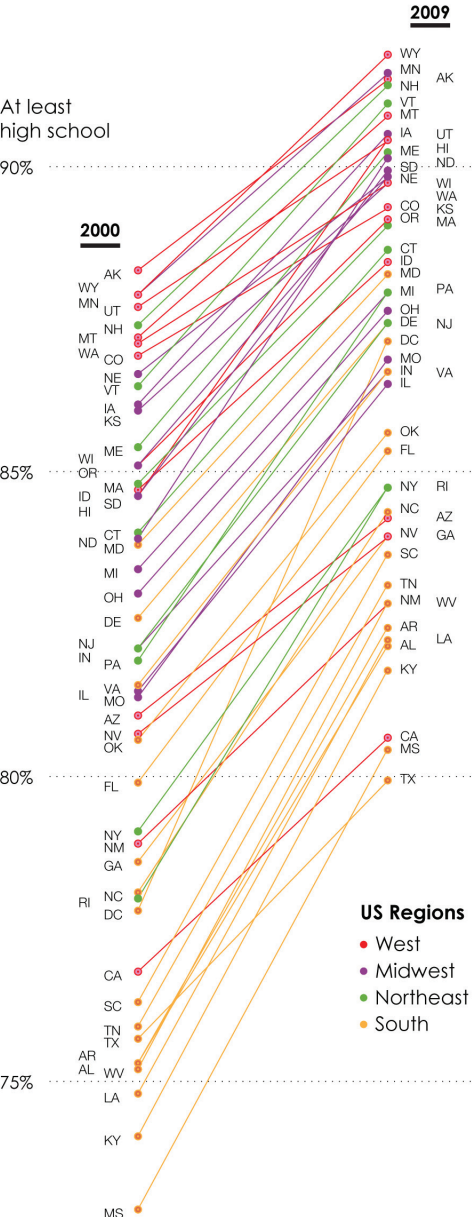
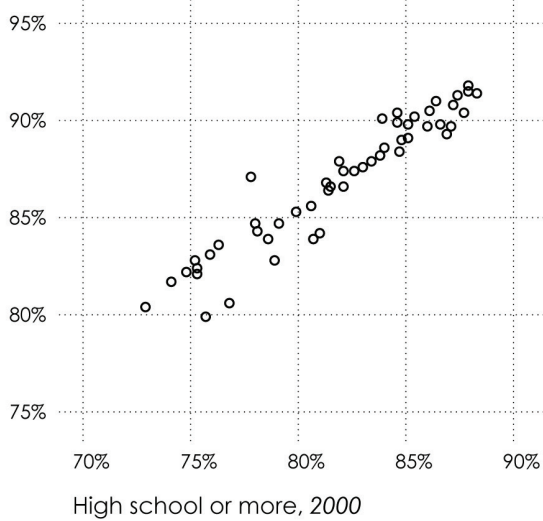


FIGURE 3-29 Using position and direction

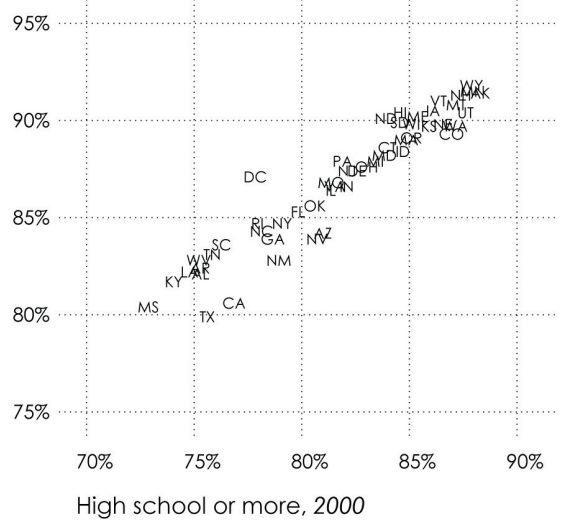
## Position

High school  
or more,  
in 2009



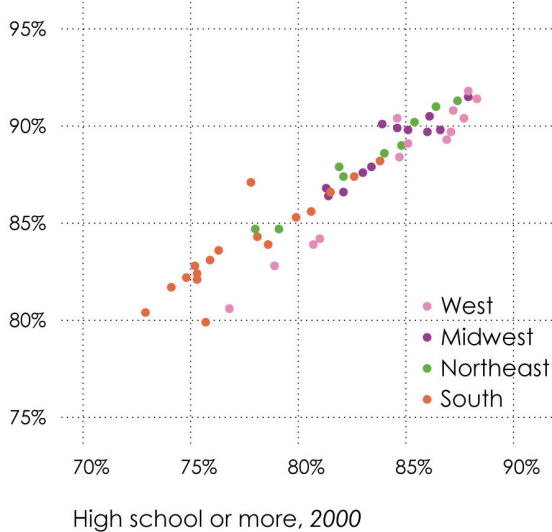
## Position + Symbols

High school  
or more,  
in 2009



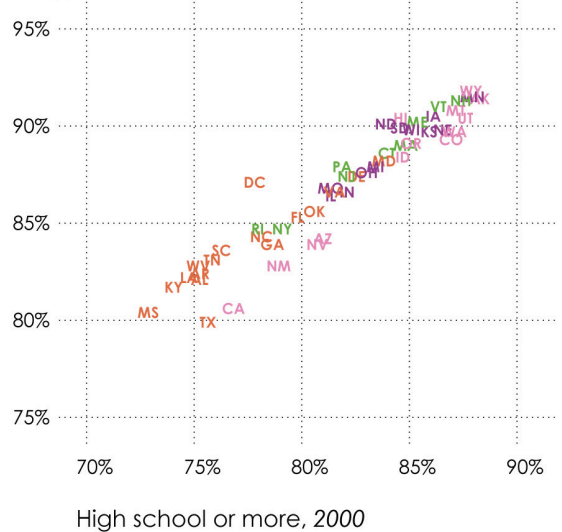
## Position + Color

High school  
or more,  
in 2009



## Position + Symbols + Color

High school  
or more,  
in 2009



**FIGURE 3-30** Position, symbols, and colors in scatterplots

High school attainment in 2000 is plotted on the horizontal axis, and attainment in 2009 is on the vertical. There is an obvious upward trend, and you can spot Washington, DC sticking out somewhat, indicating the higher rate of improvement (and probably difference in demographics). You can also see Texas and California lagging around the bottom-left corner. As shown in previous charts—and I’m sure you’re getting the hang of it now—you can incorporate other visual cues such as color, symbols, or both to provide additional dimensions of information.

**Note:** If you want to annoy cartographers, you can also call choropleth maps heat maps, as they are often referred to. The heat map was created to visualize 2-D data, and choropleth maps are the geographic equivalent. I personally keep the terms and methods separate.

Remember this is geographic data, so you must map it, right? (Actually, just because location is attached to your data, which seems like almost always these days, a map is not always the most useful view, which is discussed in the next chapter.) Figure 3-31 shows a handful of maps with states colored using varying scales and metrics, which are called choropleth maps.

Note that although each map uses the same method, the choice of scale can change the map’s focus and message. For example, the map on the top left uses a quartile scale, which means the states were split into four even groups based on a metric. In this case, the metric is the percentage of people with a bachelor’s degree in 2009. This makes a map with colors that are evenly distributed.

However, the map that shows the same data on a linear scale, with just three shades of green, shows darker shades in the Midwest and Northeast regions. Compare this with the quartile map, and you still get the lighter areas in the South, but the rest of the map tells a different story. Likewise, you can further abstract the data by coloring states by whether they are below or above the average (top right) or whether percentages increased or decreased (bottom right).

As shown in Figure 3-32, you can also show several maps at once to see how something has changed geographically over time. Since you’ve looked the data from several perspectives already, you know that a high value in 2000 generally means a higher value in 2009, because the states improved at similar rates.

You see about the same thing when you compare 1990 to 2000. In 1990, you see a more lightly colored map, where several states showed 15 percent or less of people 25 years or older with a bachelor’s degree. Only Wyoming, which had the highest percentage in 2009, shows a percentage higher than 25 percent. As you move left to right, the map gets darker, like you’d expect.

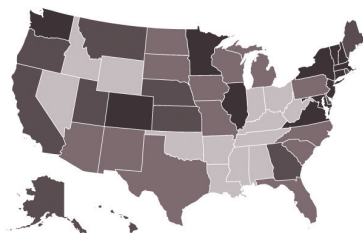


## Varying scales

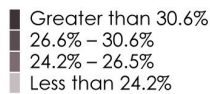
Choice of scale can shift focus and present a different message. The below maps represent how a single dataset can easily change based on this choice.

### Quartiles

Breaks decided by splitting into four equally-sized groups

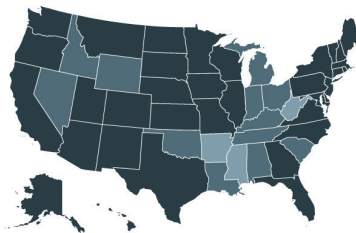


% with at least Bachelor's in 2009

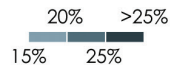


### Linear

Scale incremented evenly over range

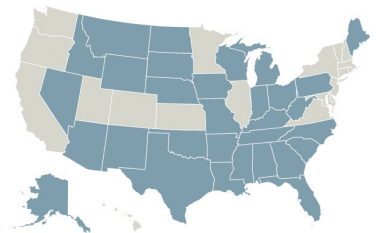


% with at least Bachelor's in 2009



### Numeric category

Create category based on a metric in data

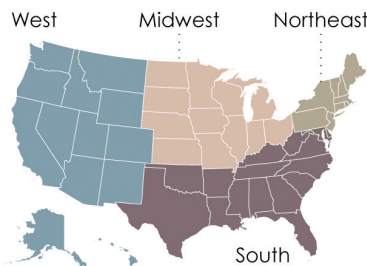


2009 US Avg. of 27.9%

Below avg. Above avg.

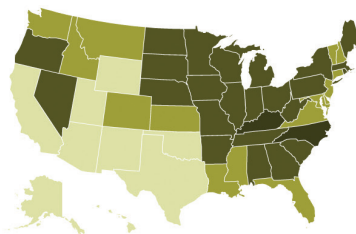
### Categorical

Groups based on metadata, such as region

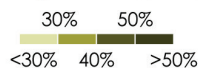


### Difference

A linear scale, but based on percent change between years

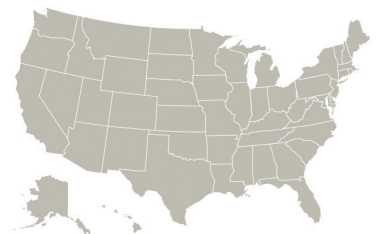


% change from 1990 to 2009



### Categorical difference

Simple split based on increase or decrease (Good news: all increase in this example)



Change from 1990 to 2009



FIGURE 3-31 Choropleth maps

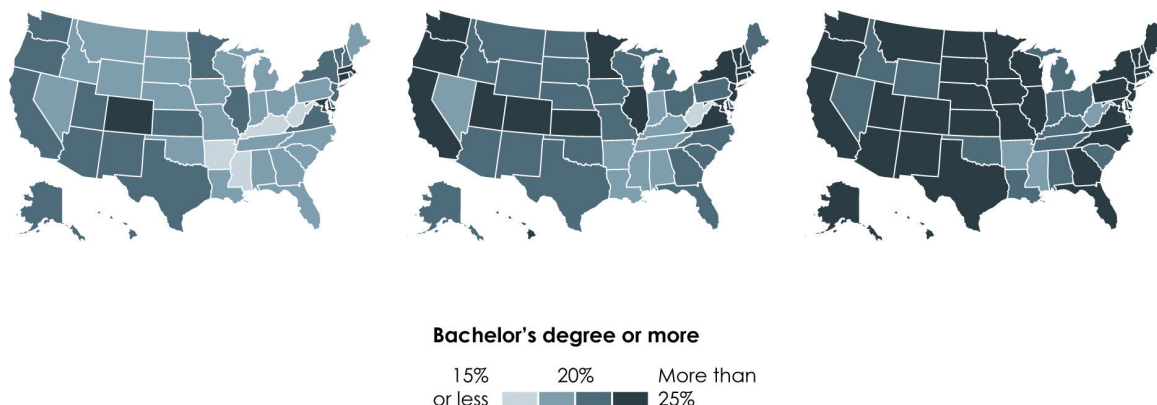


## Geographic coordinates + Time scale + Color

**In 1990**, about 20 percent of 25+ year olds had at least a bachelor's degree.

**In 2000**, the percentage was up to 24.

**In 2009**, the US average was up to 29 percent.



**FIGURE 3-32** *Maps over time*

## WRAPPING UP

At its core, visualization is an abstraction process to map data to color and geometry. This is easy to do from a technical perspective. You can easily draw and color shapes with a pencil and paper. The challenge is to figure out what shapes and colors work best, where to put them, and how to size them.

To make the jump from data to visualization, you must know your ingredients. A skilled chef doesn't just blindly throw ingredients in a pot, turn the stove on high, and hope for the best. Instead, the chef gets to know how each ingredient works together, which ones don't get along, and how long and at what temperature to cook these ingredients.

With visualization, visual cues, coordinate systems, scales, and context are your ingredients. Visual cues are the main thing that people see, and the coordinate system and scale provide structure and a sense of space. Context breathes life into the data and makes it understandable, relatable, and worth looking at.

Get to know how the components work, play with them, and get other people to look at your results and see what information they extract.

Don't forget the main component of every visualization, though. Without data, you have nothing to visualize. Likewise, if you have data with little substance, you get visualization with little substance. However, when you do get data that offers a high number of dimensions or is granular enough to see the interesting details, you still must know what to look for.

The challenge of more data is that you have more visualization options, and many of those options will be poor ones. To filter out the bad and find the worthwhile options—to get to visualization that means something—you must get to know your data. Now on to exploration.