

Data Visualization

Content:

- What is Data Visualization
- Why Data Visualization matters
- Information of Data Visualization in R
- Tools of Data Visualization in R
- Example: How to make pie-charts in R



Introduction:

“Visuals are a powerful way to convey message present information, and persuade audiences (1).”

Do you agree or disagree? Before you answer my question, there are some facts and examples.

- Human's eyes see images in less than **1/10** of a second.
- In the brain, the process of visuals is **60,000** faster than the process of text.
- The **most** information transmitted to the brain is visual.
- Human's eyes review **36,000** visual information per hour.

The above examples show the importance and the power of the visuals. The next question is how should we apply visuals in our project? The answer is **data visualization**. At the begin of this class, we have already learned that visualization is an important part of data analysis. In this post, I will explain data visualization in details. “**what is data visualization?**” and “**why is data visualization important?**” This post will answer two questions. It will also contain the information of data visualization in R, especially the tools. As example, I will show how to make pie-chat in R.

What is Data Visualization?

“This birth of statistical thinking was also accompanied by a rise in visual thinking: diagrams were used to illustrate mathematical proofs and functions; nomograms were developed to aid calculations; various graphic forms were invented to make the properties of empirical numbers— their trends, tendencies, and distributions— more easily communicated, or accessible to visual inspection (4).”

“The main goal of data visualization is its ability to visualizedata, communicating information clearly and effectively. (5)”

Basically, data visualization is a method/ presentation of data in a pictorial format, such as statistical graphics, plots and information graphics. The main goal of data visualization is to present information *clearly* and *efficiently*. Moreover, “*It involves the creation and study of the visual representation of data (2).*” It is both an art and a science. Data visualization helps people to see analytics in visual, which is easier to grasp difficult identify and confusing concepts in big data. Most of data visualization are created for human consumption. It is interesting that human can distinguish those differences in line length, shape, color and direction without any sign post. Data visualization gives senses in intuitive way.

Why Data Visualization matters?

“In other words, we see – and use and use – more of our data more of our data and process it with visual pattern recognition as the basis of o and process it with visual pattern recognition as the basis of our interpretation (8).”

Visual is a powerful way to present information. Especially, human's brains process visual information in the fastest speed. Using plots and graphs to visualize large data is intuitive and efficient. “Data visualization is going to change the way our analysts work with data. They're going to be expected to respond to issues more rapidly. And they'll need to be able to dig for more insights – look at data differently, more

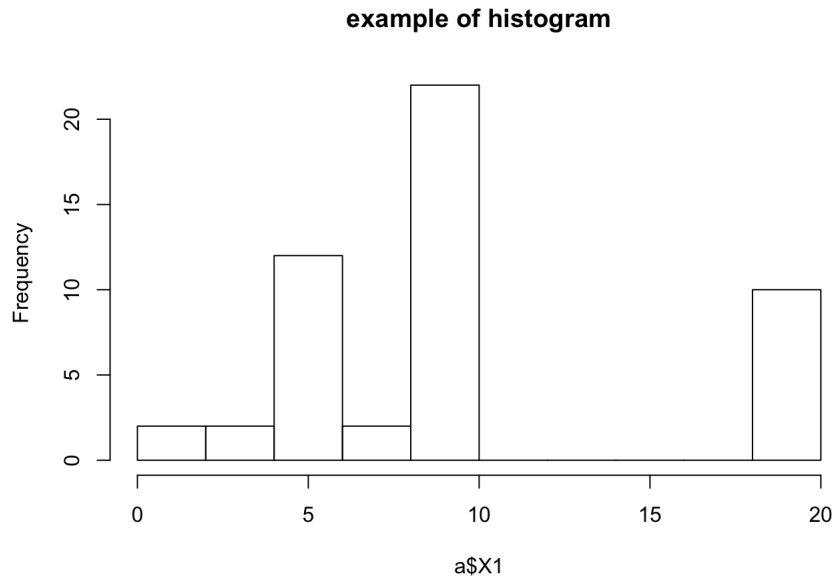
imaginatively. Data visualization will promote that creative data exploration. (Simon Samuel)"

Information of Data Visualization in R:

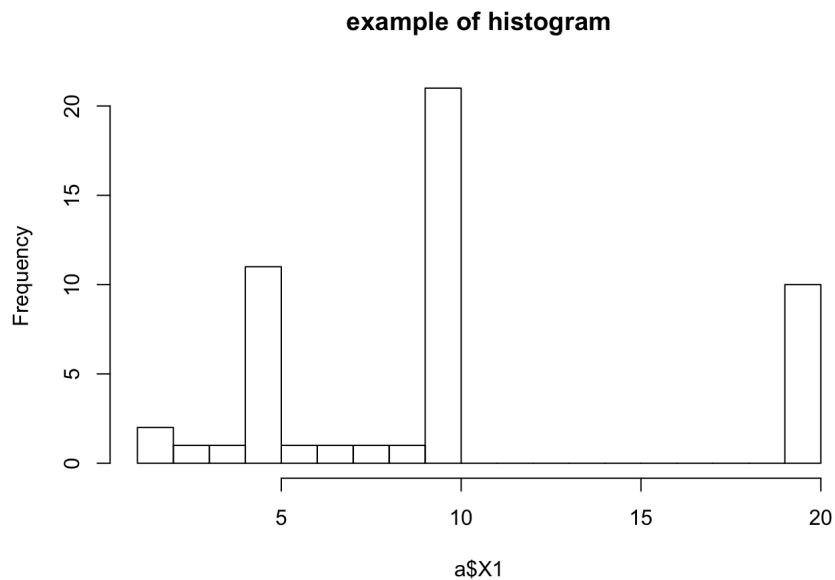
- The very basic visualization in R:

- Histogram: it is a plot that breaks the data into bins, which shows the frequency distribution of these bins. Example:

```
# example of histogram
a <- data.frame("1" = c(1:10, rep(10, 20), rep(20, 10), rep(5, 10)), "2" = 1:50)
hist(a$X1, breaks = 10, main = "example of histogram")
```



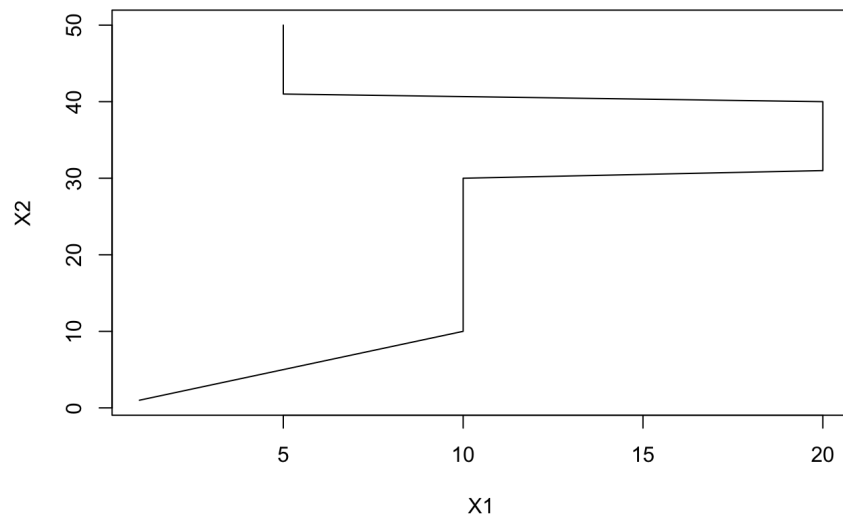
```
# different breaks
hist(a$X1, breaks = 20, main = "example of histogram")
```



- Line Chart: it is a plot that shows the trend spread over a time period. It is also a good presentation in comparing relative changes. Example:

```
# example of line
plot(a, type = "l", main = "example of line")
```

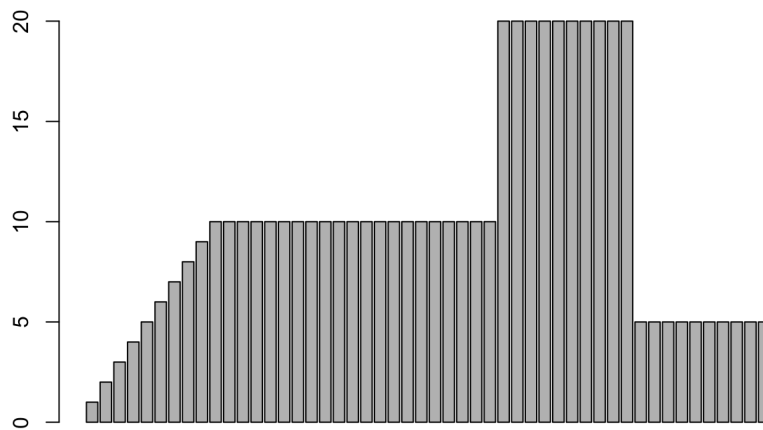
example of line



- Bar Chart: Bar plots are also suitable for showing comparisons like line chart. Moreover, it shows the cumulative totals across several groups. Example:

```
#example of bar chart  
barplot(a[X1, main = "example of bar chart"])
```

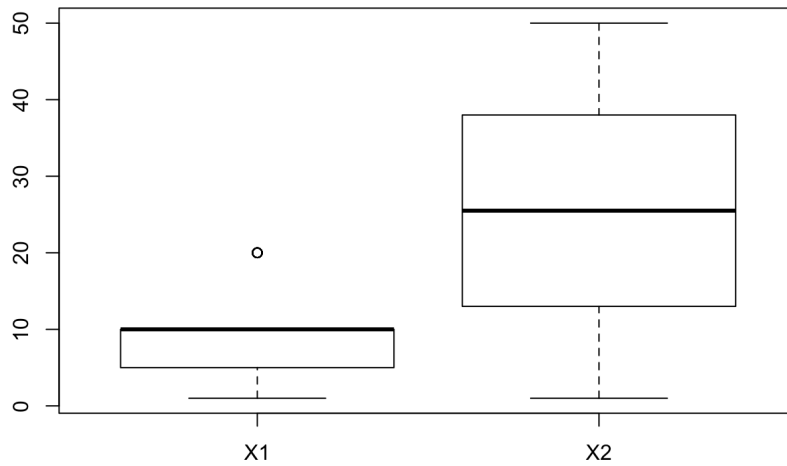
example of bar chart



- Box Plot: It is thus useful for visualizing the spread of the data and deriving inferences accordingly (3). There are five statistically significant numbers. They are the minimum, the 25th percentile, the median, the 75th percentile and the maximum. Example:

```
#example of boxplot  
boxplot(a, main = "example of boxplot")
```

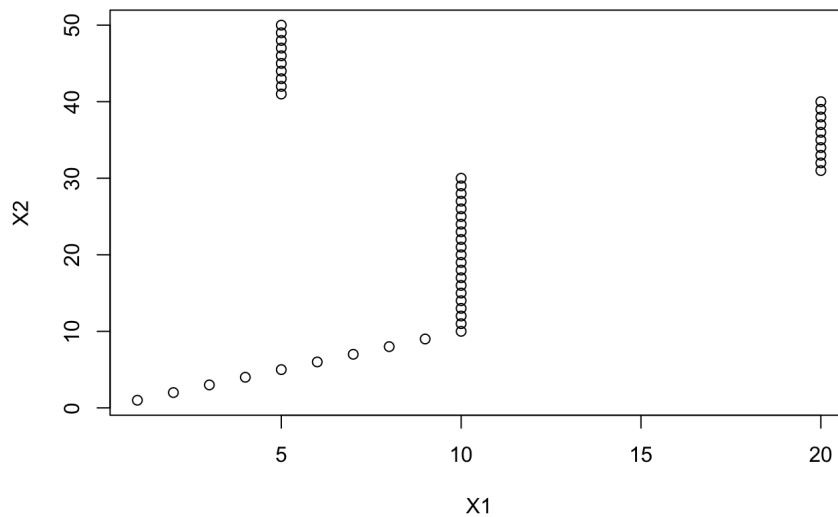
example of boxplot



- Scatter plot: it helps in visualizing data intuitively and for simple data inspection. It shows the distribution of data. Example:

```
# example of plot  
plot(a, main = "example of scatter plot")
```

example of scatter plot



- The adventure visualization in R:



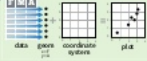
- **ggplot2:** we learn this in class, simply, it is a plotting system for R, which based on the grammar of graphics. It cares about many details and tries to provide a powerful model of graphics.

Data Visualization with ggplot2 Cheat Sheet

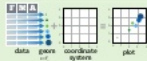


Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data set**, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x and y locations**.



Build a graph with **ggplot()** or **qplot()**

ggplot(data = mpg, aes(x = cty, y = hwy))
Begins a plot that you finish by adding layers to. No defaults, but provides more control than qplot().

geom
Add layers to a plot with a **geom_*** or **stat_*** function. Each provides a geom, a set of aesthetic mappings, and a default stat and position adjustment.

ggplot(mapping = aes(x = hwy, color = cyl), data = mpg, geom = "point")
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

last_plot()
Returns the last plot.

ggsave("plot.png", width = 5, height = 5)
Saves last plot as 5 x 5 file named "plot.png" in working directory. Matches file type to file extension.

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 999-498-1212 • rstudio.com

Geoms - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

One Variable

Continuous

a = **geom_area**(stat = "bin")
x, y, alpha, color, fill, linetype, size
b = **geom_area**(stat = "density", stat = "bin")
x, y, alpha, color, fill, linetype, size, weight
b = **geom_density**(kernel = "gaussian")
x, y, alpha, color, fill, linetype, size, weight
b = **geom_density2d**(aes(y = .density))
x, y, alpha, color, fill
a = **geom_freqpoly**()
x, y, alpha, color, linetype, size
b = **geom_freqpoly**(aes(y = .density))
a = **geom_histogram**(binwidth = 5)
x, y, alpha, color, fill, linetype, size, weight
b = **geom_histogram**(aes(y = .density))
Discrete
b = **geom_bar**()
x, y, alpha, color, fill, linetype, size, weight

Graphical Primitives

map = **map_data**("state")
c = **ggplot(map, aes(long, lat))**
d = **geom_polygon**(aes(group = group))
x, y, alpha, color, fill, linetype, size
d = **ggplot(economics, aes(date, unemployment))**
d = **geom_path**(inorder = "true", linetype = "solid", linewidth = 1)
x, y, alpha, color, linetype, size
d = **geom_ribbon**(aes(ymin = unemployment - 90, ymax = unemployment + 90))
x, y, alpha, color, fill, linetype, size
e = **ggplot(seals, aes(x = long, y = lat))**
e = **geom_segment**(aes(xend = long + delta_long, yend = lat + delta_lat))
x, y, alpha, color, linetype, size
e = **geom_rect**(aes(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat))
x, y, alpha, color, fill, linetype, size

Two Variables

Continuous X, Continuous Y

f = **geom_blank**()
(useful for expanding limits)
f = **geom_jitter**()
x, y, alpha, color, fill, shape, size
f = **geom_point**()
x, y, alpha, color, fill, shape, size
f = **geom_quantile**()
x, y, alpha, color, linetype, size, weight
f = **geom_rug**(sides = "bl")
alpha, color, linetype, size
f = **geom_smooth**(method = "lm")
x, y, alpha, color, fill, linetype, size, weight
f = **geom_text**(aes(label = cty))
x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

Discrete X, Continuous Y

g = **ggplot(mpg, aes(class, hwy))**
g = **geom_bar**(stat = "identity")
x, y, alpha, color, fill, linetype, size, weight
g = **geom_boxplot**()
lower, middle, upper, x, y, alpha, color, fill, linetype, shape, size, weight
g = **geom_dotplot**(binaxis = "y", stackdir = "center")
x, y, alpha, color, fill
g = **geom_violin**(scale = "area")
x, y, alpha, color, fill, linetype, size, weight

Discrete X, Discrete Y

h = **ggplot(diamonds, aes(cut, color))**
h = **geom_jitter**()
x, y, alpha, color, fill, shape, size

Three Variables

m = **geom_raster**(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)
x, y, alpha, fill (slow)
m = **geom_tile**(aes(fill = z))
x, y, z, alpha, color, linetype, size, weight

Continuous Bivariate Distribution

i = **ggplot(movies, aes(year, rating))**
i = **geom_bin2d**(binwidth = c(5, 0.5))
x, y, xmin, xmax, ymin, ymax, alpha, color, fill, linetype, size, weight
i = **geom_density2d**()
x, y, alpha, color, linetype, size
i = **geom_hex**()
x, y, alpha, color, fill, size

Continuous Function

j = **ggplot(economics, aes(date, unemploy))**
j = **geom_area**()
x, y, alpha, color, fill, linetype, size
j = **geom_line**()
x, y, alpha, color, linetype, size
j = **geom_step**(direction = "hv")
x, y, alpha, color, linetype, size

Visualizing error

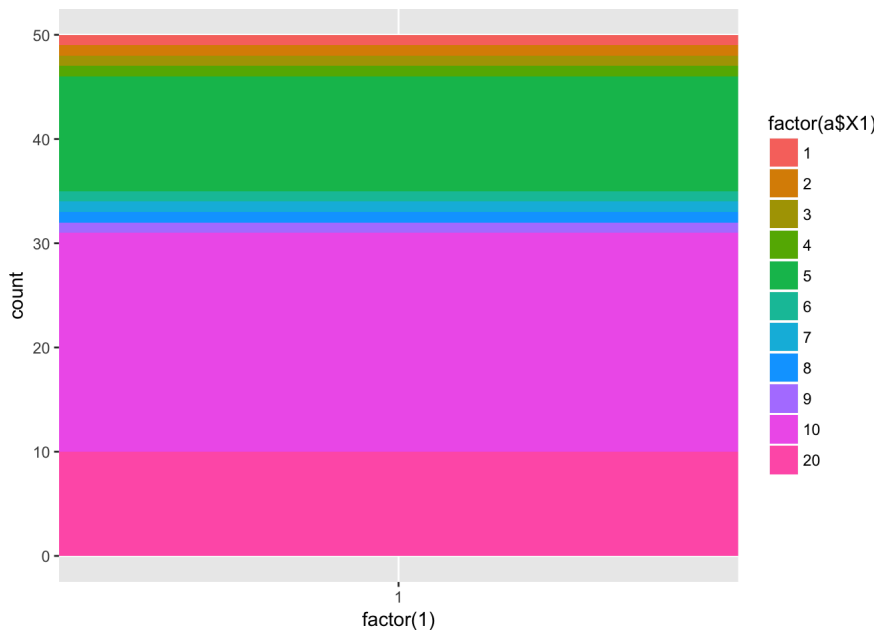
df = **data.frame**(group = c("A", "B"), fit = 4.5, se = 1.2)
k = **ggplot(df, aes(group, fit, ymin = fit - se, ymax = fit + se))**
k = **geom_crossbar**(latten = 2)
x, y, xmin, xmax, alpha, color, fill, linetype, size
k = **geom_errorbar**()
x, y, xmin, xmax, alpha, color, linetype, size, width (also **geom_errorbarh**)
k = **geom_linerange**()
x, y, xmin, xmax, alpha, color, linetype, size
k = **geom_pointrange**()
x, y, xmin, xmax, alpha, color, fill, linetype, shape, size

Maps

data = **data.frame**(murder = USArrests\$Murder, state = tolowerrownames(USArrests))
map = **map_data**("state")
i = **ggplot(data, aes(fill = murder))**
i = **geom_map**(aes(long, lat = state, map = map))
e = **expand_limits**(map_id(alpha, color, fill, linetype, size)

For example, I will show how to use ggplot2 to visualization a pie-chart.

```
# example of pie-chart in ggplot2
library(ggplot2)
# first create a bar plot
bar <- ggplot(a, aes(x = factor(1), fill = factor(a$X1))) + geom_bar(width = 1)
bar
```



```
# second transfer this bar plot into polar chart by adding coord_polar
pie <- bar + coord_polar(theta = "y")
pie
```


2. <https://cran.r-project.org/web/packages/hexbin/hexbin.pdf>
3. "Data visualization: A wise investment in your big data future." SAS, www.sas.com/en_us/insights/articles/analytics/data-visualization-a-wise-investment-in-your-big-data-future.html.
4. <http://www.math.yorku.ca/SCS/Gallery/milestone/milestone.pdf>
5. Friedman, Vitaly. "Data Visualization and Infographics." Smashing Magazine, 14 Jan. 2008, www.smashingmagazine.com/2008/01/monday-inspiration-data-visualization-and-infographics/.
6. techATstate. "Tech@State: Data Visualization - Keynote by Dr Edward Tufte." YouTube, YouTube, 7 Aug. 2013, www.youtube.com/watch?v=g9Y4SxgfGCg.
7. https://www.sfu.ca/gis/geog_x55/web355/icons/11_lec_vweb.pdf
8. Quoted in M. Wood. 1994. "The Traditional Map." in Visualization in Geographic Information Systems, Edited by Hilary Hearnshaw & David Unwin. John Wiley & Sons.
9. https://www.sas.com/en_us/insights/big-data/data-visualization/_jcr_content/socialShareImage.img.png
10. http://www.dmeforpeace.org/sites/default/files/images/081114_FeatureBlog_0.png

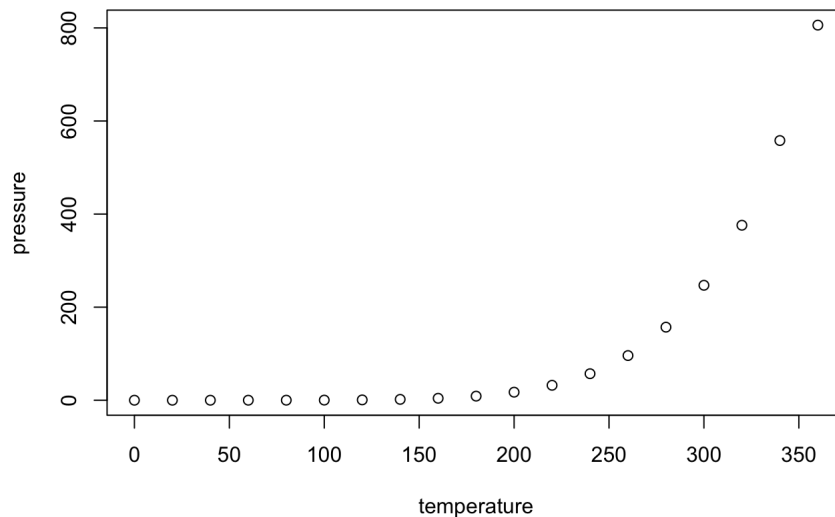
R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

Including Plots

You can also embed plots, for example:



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.