

零基础学数据可视化

一、导论

毛益民 副教授

浙江工商大学公共管理学院

2020-02-21

为何学数据可视化?



图表的作用

- 真实、准确、全面地展示数据；
- 以较小的空间承载较多的信息；
- 揭示数据的本质、关系、规律。

--

可视化的终极目标是洞悉蕴含在数据中的现象和规律。

这包括多重含义：发现、决策、解释、分析、探索和学习.

正所谓：“一图胜千言”

Data Scientist The Sexy Job



October 2012 Issue

DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

When Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know anyone. So you just stand in the corner sipping your drink—and you probably leave early."

- See also an old article by NYT (2009): **For Today's Graduate, Just One Word: Statistics**
- And another famous McKinsey 2011 Report: **Big data: The next frontier for innovation, competition, and productivity**

为何采用R做数据分析？

🔔 R是什么？



🐾 R能做什么？

小调查

- 1.有处理过数据吗？量比较大，内容比较复杂那种？
- 2.你们平时用什么软件处理数据？

绘图软件比较

LOGO	名称	开源	付费	技能要求	官方网站
	Excel	否	是	界面操作	https://support.office.com/en-GB/Excel
	Origin	否	是	界面操作	http://originlab.com/
	SigmPlot	否	是	界面操作	https://systatsoftware.com/products/sigmaplot/
	GraphPad Prism	否	是	界面操作	http://www.graphpad.com/
	MATLAB	否	是	编程	https://www.mathworks.com/products/matlab.html
	Python	是	否	编程	https://www.python.org/
	R	是	否	编程	https://www.r-project.org/

R是什么？

R is a language and environment for statistical computing and graphics. It includes

- an effective data handling and storage facility,
- a suite of operators for calculations on arrays, in particular matrices,
- a large, coherent, integrated collection of intermediate tools for data analysis, graphical facilities for data analysis and display either on-screen or on hardcopy, and
- a well-developed, simple and effective programming language which includes conditionals, loops, user-defined recursive functions and input and output facilities.

The term “environment” is intended to characterize it as a fully planned and coherent system, rather than an incremental accretion of very specific and inflexible tools, as is frequently the case with other data analysis software.

Rstudio IDE

The screenshot shows the RStudio IDE interface with the following components:

- Script Editor:** Contains R code for generating a diamond pricing plot. The code includes loading ggplot2, summarizing the diamonds dataset, creating a mean carat size variable, and plotting carat vs price by clarity.
- Workspace:** Shows the diamonds dataset (53,940 observations) and a ggplot object named p.
- Console:** Displays the summary statistics for the diamonds dataset, including the minimum, first quartile, median, mean, third quartile, and maximum values for carat, price, and depth.
- Plots:** A scatter plot titled "Diamond Pricing" showing Price (Y-axis, 0 to 15,000) versus Carat (X-axis, 0.0 to 3.5). The data points are colored by Clarity, with categories I1, SI2, SI1, VS2, VS1, VVS2, VVS1, and IF represented by different colors.

R能做什么？

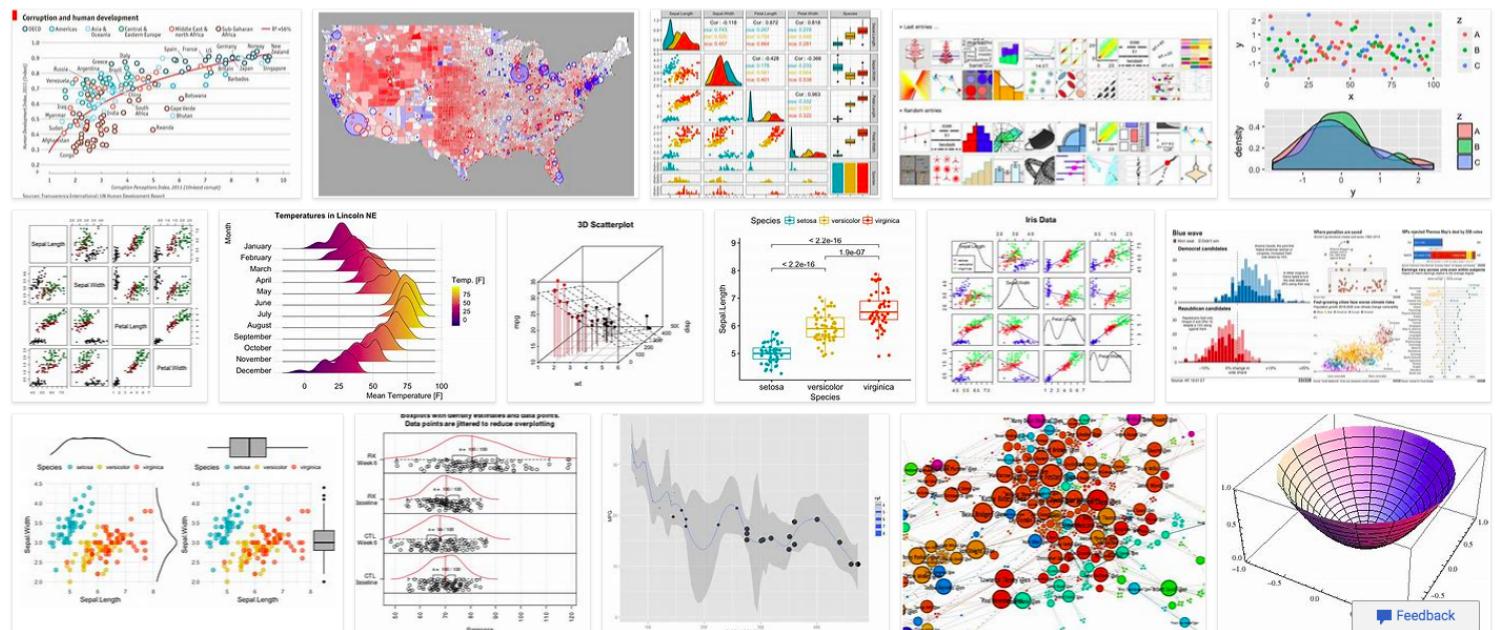
这个问题问得好？

答案是：

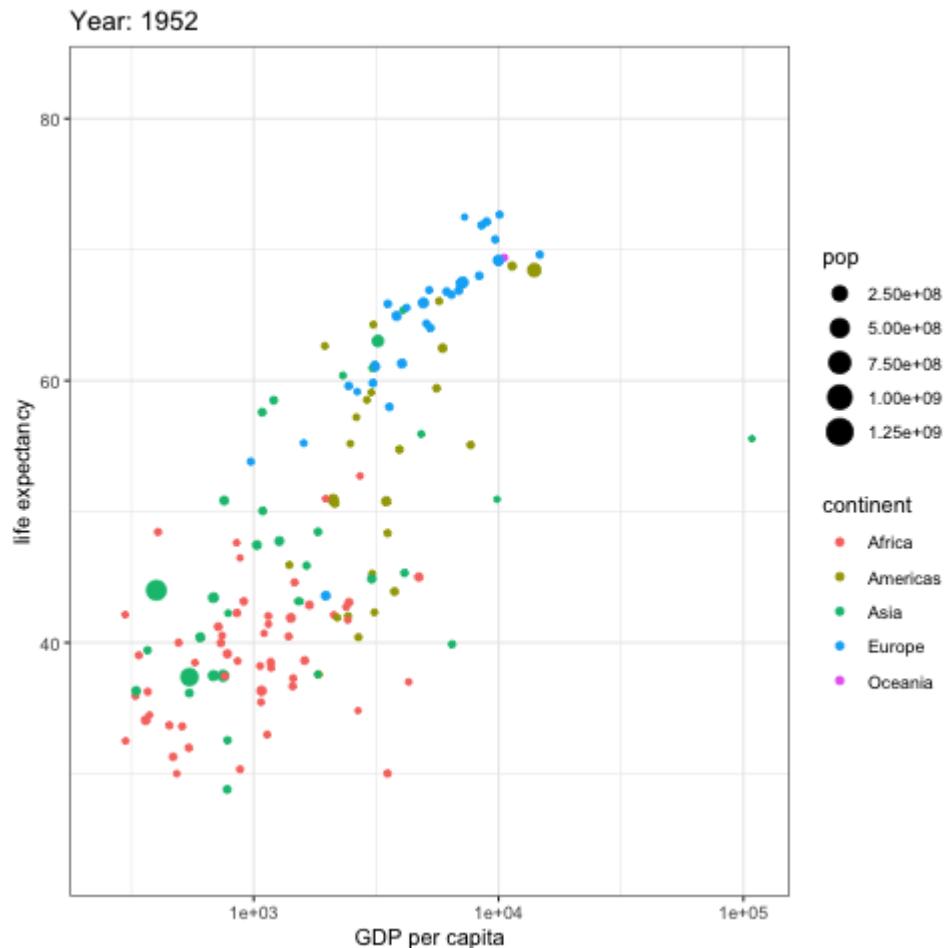
在数据分析领域，无所不能

就算现在不能，以后也一定能

用R绘制的图形



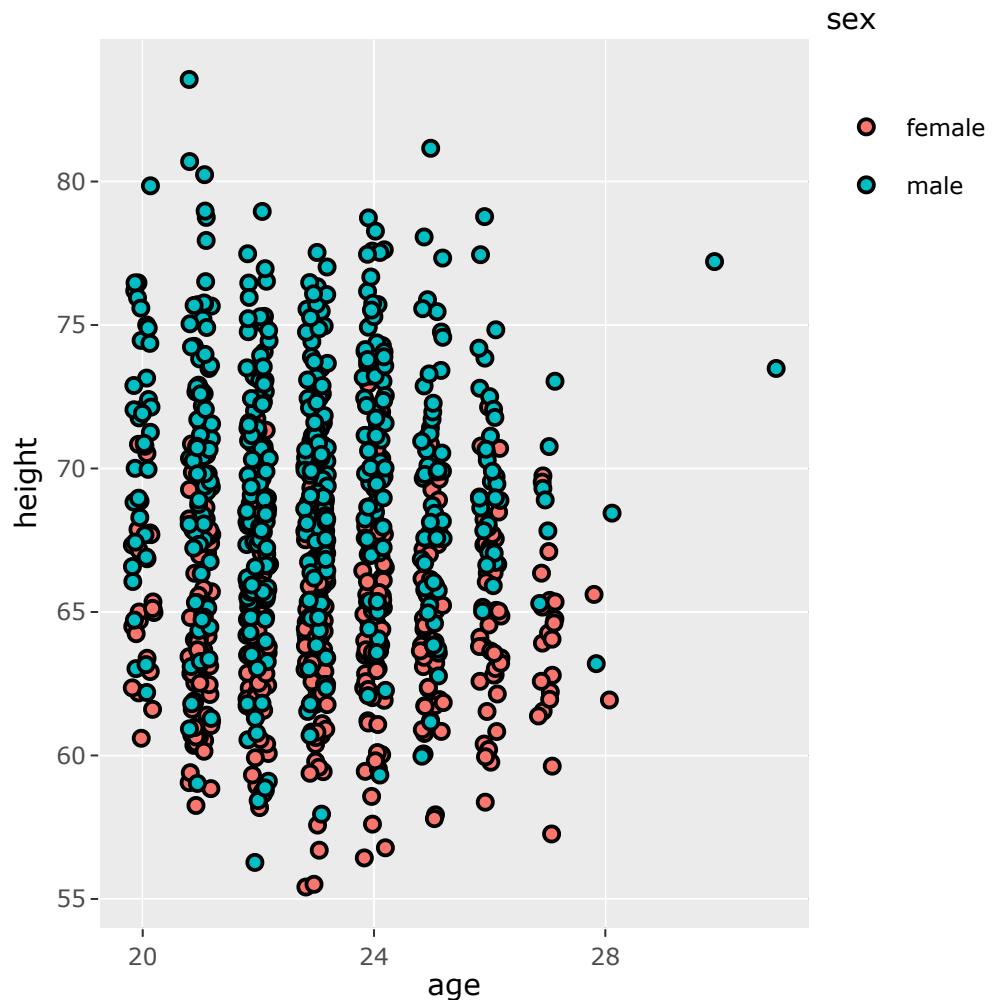
用R绘制的图形



用R绘制的图形

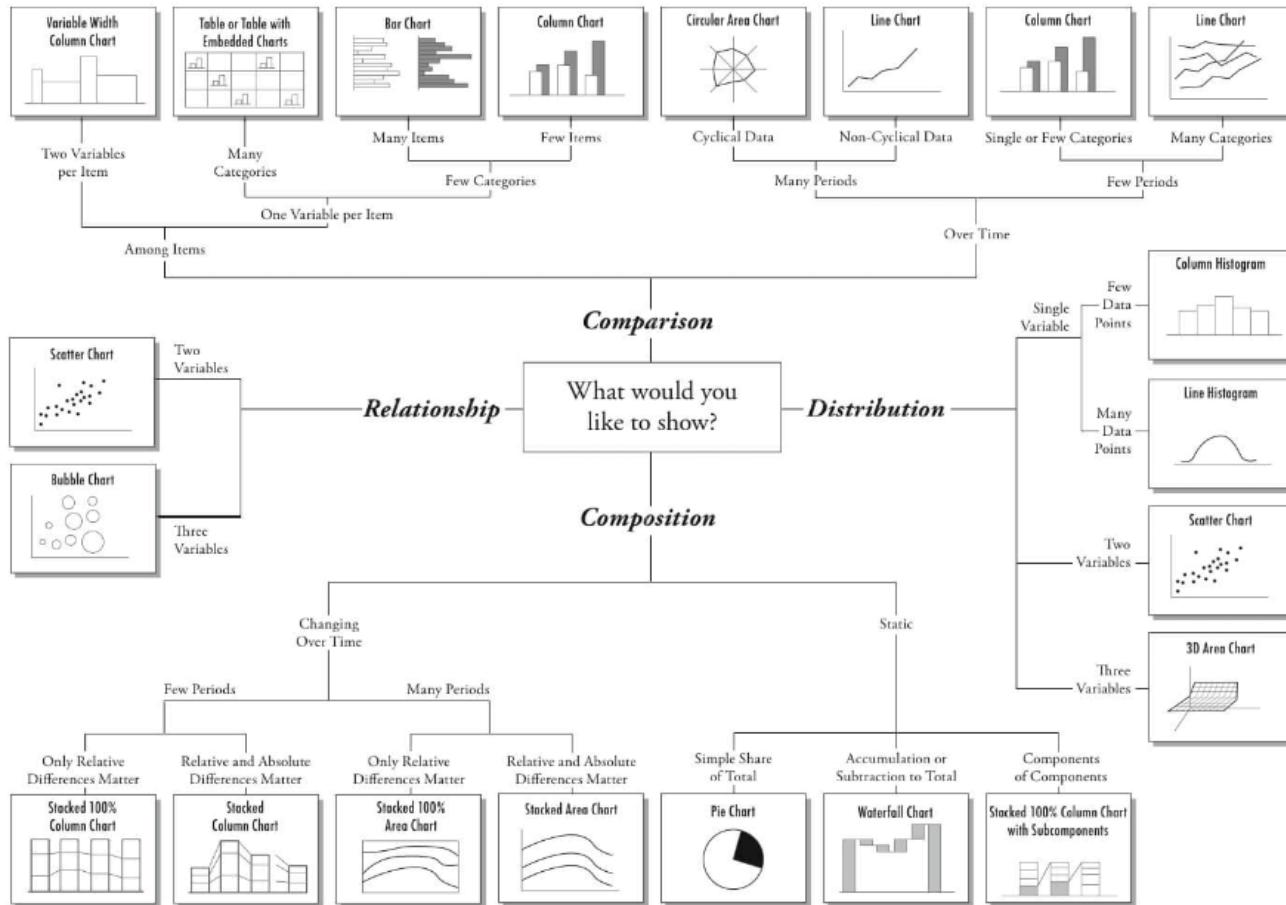
用R绘制的图形

用R绘制的图形



各种图形

Chart Suggestions—A Thought-Starter



R 代码

```
#做个计算题  
5+5+7+9+10
```

```
# [1] 36
```

```
# 生成十个数字  
rnorm(10,mean = 0,sd = 1)
```

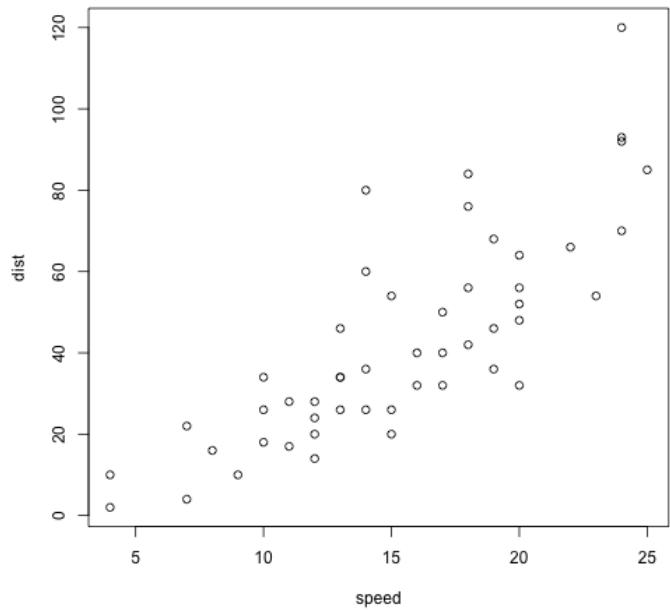
```
# [1] 1.30709064 0.03676991 -0.85136309 0.19612334 -0.02014146  
# [6] 1.42881203 0.86934421 -1.42369295 0.11197427 0.55488916
```

```
# 处理文字  
dojutsu = c('地爆天星', '天照', '加具土命', '神威', '須佐能乎', '無限月読')  
grep('天', dojutsu, value = TRUE)
```

```
# [1] "地爆天星" "天照"
```

R 基本绘图

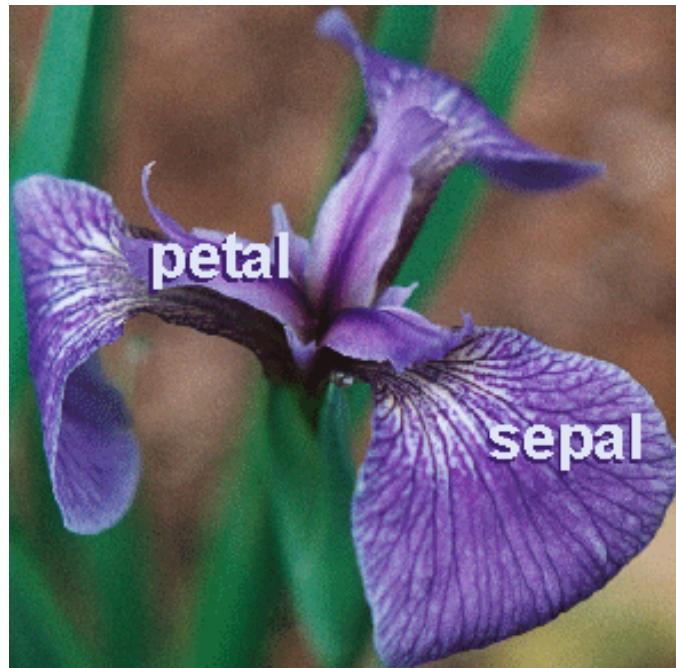
```
plot(cars)
```



R绘图

Iris Data Set （鸢尾属植物数据集）

Iris Data Set（鸢尾属植物数据集）首次出现在著名的英国统计学家和生物学家Ronald Fisher 1936年的论文《The use of multiple measurements in taxonomic problems》。



Iris Data Set （鸢尾属植物数据集）

在这个数据集中，包括了三类不同的鸢尾属植物：**Iris Setosa**，**Iris Versicolour**，**Iris Virginica**。每类收集了50个样本，因此这个数据集一共包含了150个样本。

该数据集测量了所有150个样本的4个特征，分别是：

- sepal length (花萼长度)
- sepal width (花萼宽度)
- petal length (花瓣长度)
- petal width (花瓣宽度)

以上四个特征的单位都是厘米 (cm)。通常使用 m 表示样本量的大小， n 表示每个样本所具有的特征数。因此在该数据集中， $m=150, n=4$

表格显示

Show 6 entries

Search:

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa

Showing 1 to 6 of 40 entries

Previous

1

2

3

4

5

6

7

Next

做个练习

三类不同的鸢尾属植物：

- Iris Setosa
- Iris Versicolour
- Iris Virginica

4个特征：

- sepal length (花萼长度)
- sepal width (花萼宽度)
- petal length (花瓣长度)
- petal width (花瓣宽度)

问题：请按种类分别求各特征的均值

让R来帮你完成

```
library(dplyr)
datatable(iris %>%
  group_by(Species) %>%
  summarise_all(mean))
```

Show 10 entries

Search:

	Species	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
1	setosa	5.006	3.428	1.462	0.246
2	versicolor	5.936	2.77	4.26	1.326
3	virginica	6.588	2.974	5.552	2.026

Showing 1 to 3 of 3 entries

Previous

1

Next

数据分析的基本流程

谢谢聆听!

毛益民

浙江工商大学公共管理学院

下次请带上电脑，我们一起感受R的乐趣！

各位老师、各位同学，大家好！

首先感谢浙江工商大学公共管理学院给我这次机会，让我能够在这里跟大家分享一下我使用R语言的一些心得。

我叫毛益民，是浙江工商大学公共管理学院的一名硕士研究生。我的研究方向是数据挖掘，主要利用R语言进行数据分析和模型构建。

在开始分享之前，我想先给大家简单介绍一下R语言。R语言是一种统计编程语言，广泛应用于数据处理、统计分析、机器学习等领域。它的特点是语法清晰、功能强大、易于学习和使用。

那么，我们今天要聊的是什么内容呢？

首先，我会简要介绍R语言的基本语法和常用函数，让大家对R有一个初步的了解。

其次，我会通过一些具体的例子，展示如何使用R语言进行数据分析和模型构建。

最后，我会分享一些我在使用R语言过程中遇到的问题和解决方法，希望大家能够从中受益。

好了，接下来就让我们一起进入今天的主题吧！