

A Beginner's Guide to Learn R Programming

许忠平



R version 3.3.3 (2017-03-06) -- "Another Canoe"
Copyright (C) 2017 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

```
> library(ggplot2)
> ggplot(l11z, aes(Sepal.Length, Sepal.Width)) + geom_point()
> |
```

Sepal.Width

Sepal.Length

C:/D-Datos/Göttingen/Papers/LIDAR variables selection Edu/Box-Cox - RStudio

```
# biomass calculation per tree
201 # biomass calculation per tree
202 kalimantansw.brown<-brown.m01st.d(kalimantansdbh)
203 kalimantansw.yamakura<-yamakura.stem(kalimantansdbh, kalimantansh)+yamakura.branch(yamakura.stem(k
204 kalimantansw.basuki<-basuki.mixed.d(kalimantansdbh)
205 kalimantansw.samalca<-samalca.d(kalimantansdbh)
206 kalimantansw.hashimoto<-hashimoto.d(kalimantansdbh)
207 kalimantansw.kenzo<-kenzo.d(kalimantansdbh)
208 kalimantansw.forda<-forda.d(kalimantansdbh)
209 kalimantansw.jaya<-jaya.d(kalimantansdbh)
210 kalimantansw.novita<-novita.d(kalimantansdbh)
211 kalimantansw.nugroho.d<-nugroho.d(kalimantansdbh)
212 kalimantansw.nugroho.d.h<-nugroho.d.h(kalimantansdbh)
213
214 plot(kalimantansdbh, kalimantansw.brown, col=1)
215 points(kalimantansdbh, kalimantansw.yamakura, col=2)
216 points(kalimantansdbh, kalimantansw.basuki, col=3)
217 points(kalimantansdbh, kalimantansw.hashimoto, col=5)
218 points(kalimantansdbh, kalimantansw.kenzo, col=6)
219 points(kalimantansdbh, kalimantansw.forda, col=7)
220 points(kalimantansdbh, kalimantansw.jaya, col=8)
221 points(kalimantansdbh, kalimantansw.novita, col=9)
222 points(kalimantansdbh, kalimantansw.nugroho.d, col=10)
223 points(kalimantansdbh, kalimantansw.nugroho.d.h, col=11)
224
225 legend(10,8000, c("Brown", "Yamakura", "Basuki", "Samalca", "Hashimoto", "Kenzo", "Forda", "Jaya",
226
227
228
229 # Summing all values per plot and nested plot
230 bio.plot.brown<-as.data.frame(capply(kalimantansw.brown, list(kalimantansplot_id, kalimantansubpl
231
310:1 (Untitled) R Script
```

R script

R console

Environment History

Object	Size
h1.l.trees	716 obs. of 23 variables
kal.plot	94 obs. of 18 variables
kalimantan	1993 obs. of 44 variables
l1z.plots	59 obs. of 19 variables

Biomass estimation per plot with different models

Graphical output



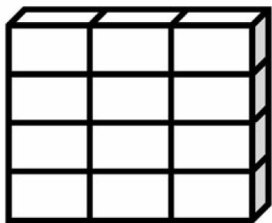
数据类型 (Data Types)

(a) 向量



```
a <- c(1, 2, 5, 3, 6, -2, 4)
b <- c("one", "two", "three")
c <- c(TRUE, TRUE)
```

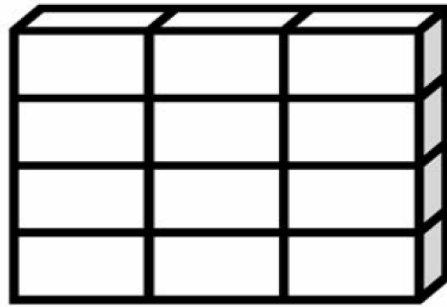
(d) 数据框



各列的模式 (modes) 可以不同

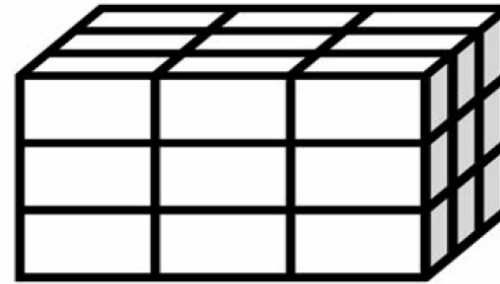
```
patientID <- c(1, 2, 3, 4)
age <- c(25, 34, 28, 52)
diabetes <- c("Type1", "Type2", "Type1", "Type1")
status <- c("Poor", "Improved", "Excellent", "Poor")
patientdata <- data.frame(patientID, age, diabetes, status)
patientdata
  patientID  age diabetes  status
1         1   25   Type1   Poor
2         2   34   Type2  Improved
3         3   28   Type1  Excellent
4         4   52   Type1   Poor
```

(b) 矩阵



```
y <- matrix(1:12, nrow=3, ncol=4)
y
  [,1] [,2] [,3] [,4]
[1,] 1   4   7  10
[2,] 2   5   8  11
[3,] 3   6   9  12
```

(c) 数组



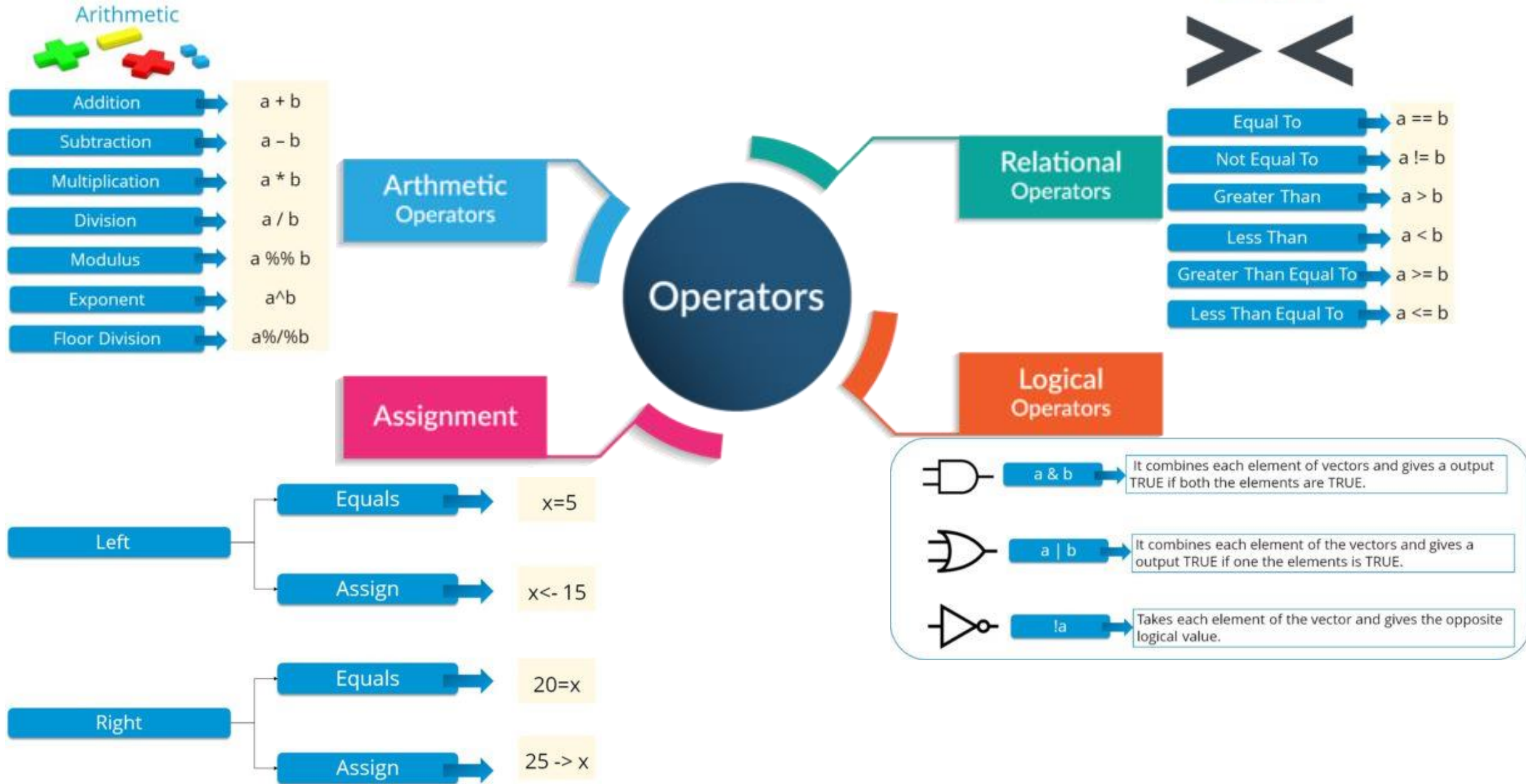
```
dim1 <- c("A1", "A2")
dim2 <- c("B1", "B2", "B3")
dim3 <- c("C1", "C2")
z <- array(1:12, c(2, 3, 2),
          dimnames=list(dim1, dim2, dim3))
z
  , , C1
     B1 B2 B3
A1 1   3   5
A2 2   4   6
  , , C2
     B1 B2 B3
A1 7   9  11
A2 8  10  12
```

(e) 列表

向量
数组
数据框
列表

```
g <- "My First List"
h <- c(25, 26, 18, 39)
j <- matrix(1:4, nrow=2)
k <- c("one", "two", "three")
mylist <- list(title=g, ages=h, j, k)
mylist
$title
[1] "My First List"
$ages
[1] 25 26 18 39
[[3]]
  [,1] [,2]
[1,] 1   3
[2,] 2   4
[[4]]
[1] "one" "two" "three"
```

数据操作 (Data Operators)



函数(Functions)

求平均: `mean(x)`

求和: `sum(x)`

开平方: `sqrt(x)`

.....



```
1 > example("mean")
2
3 mean> x <- c(0:10, 50)
4
5 mean> xm <- mean(x)
6
7 mean> c(xm, mean(x, trim = 0.10))
8 [1] 8.75 5.50
9
```

```
function_name <-function(arg_1, arg_2, ...) {
//Function body
}
```

```
1 sum_of_square <- function(x,y) {
2   x^2 + y^2
3 }
4 sum_of_squares(3,4)
5
```


R 包 (R packages)

R packages are collections of functions and data sets developed by the community.



CRAN

```
install.packages("XXX")
```



Bioconductor

```
source("https://bioconductor.org/biocLite.R")  
biocLite("XXX")
```



Github

```
install.packages("devtools")  
devtools::install_github("hadley/babynames")
```

library(XXX)

help(XXX)

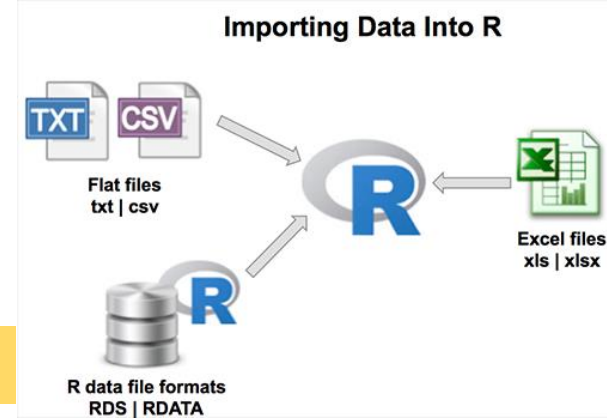
search()



读取数据 (Import from files)

R base functions:

```
1 # Command to copy&paste tables from Excel or other programs into R.
2 read.delim("clipboard", header=T)
3 -----
4 # Reads in table and assigns it to data frame, with info on column headers and field separators.
5 read.table(file="file", header=TRUE, sep="\t")
6 -----
7 # Reads a file in table format and creates a data frame from
8 read.csv("file",header = T)
9
```



.txt

.CSV

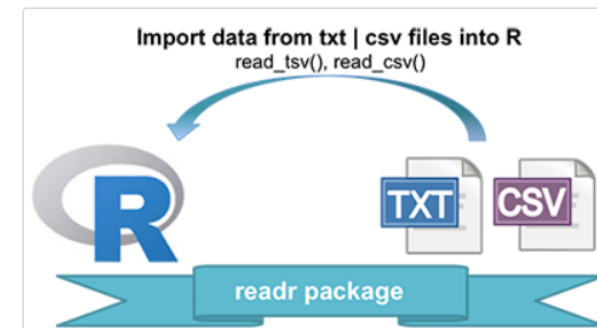
- file参数：这是必须的，可以是相对路径或者绝对路径（注意：Windows下路径要用斜杠 '/' 或者双反斜杠 '\\'）。
- 如果数据集中含有中文，直接导入很有可能不识别中文，这时加上参fileEncoding='utf-8'

R package (readr):

```
1 install.packages("readr")
2 # Loading
3 library("readr")
4 read_csv(file, col_names = TRUE)
5 read_tsv(file, col_names = TRUE)
6
```

Compared to R base functions, readr functions are:

- ❖ much faster (X10),
- ❖ have a helpful progress bar if loading is going to take a while
- ❖ all functions work exactly the same way.



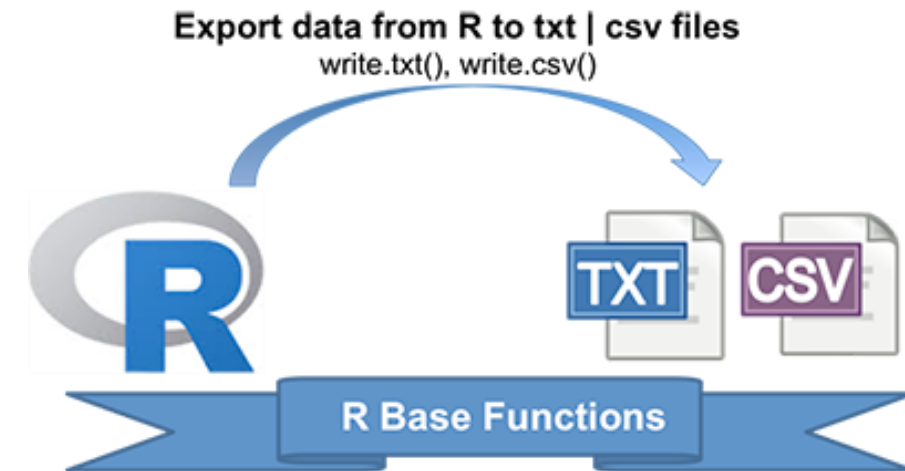
导出数据 (Export data from R)

R base functions:

```
1 data("mtcars")
2 write.table(mtcars, file = "mtcars.txt", sep = "\t",
3 |           row.names = TRUE, col.names = NA)
4 -----
5 write.csv(mtcars, file = "mtcars.csv")
6
7
```

R package (readr):

```
1 write_csv(challenge, "challenge.csv")
2
```



数据转换 (Data transformation)

清洗和整理

```
1 library(nycflights13)
2 library(tidyverse)
3 # 1.1 筛选: filter()
4 (jan1 <- filter(flights, month == 1, day == 1))
5 # 1.2 排列: arrange()
6 arrange(flights, year, month, day)
7 arrange(flights, desc(arr_delay))
8 # 1.3 选择: select()
9 select(flights, year, month, day)
10 # 1.4 变形: mutate()
11 flights_sml <- select(flights,
12   year:day,
13   ends_with("delay"),
14   distance,
15   air_time )
16 ##### 新添加的列可以用于后续计算
17 mutate(flights_sml,
18   gain = arr_delay - dep_delay,
19   hours = air_time / 60,
20   gain_per_hour = gain / hours )
21 ##### 只保留变形后的列
22 transmute(flights,
23   gain = arr_delay - dep_delay,
24   hours = air_time / 60,
25   gain_per_hour = gain / hours )
26 # 1.5 汇总: summarise()
27 summarise(flights, delay = mean(dep_delay, na.rm = TRUE))
28 # 1.6 分组: group_by()
29 by_day <- group_by(flights, year, month, day)
30 summarise(by_day, delay = mean(dep_delay, na.rm = TRUE))
31
```



筛选

按给定的逻辑判断筛选出符合要求的子数据集

排列

按给定的列名依次对行进行排序

选择

用列名作参数来选择子数据集

变形

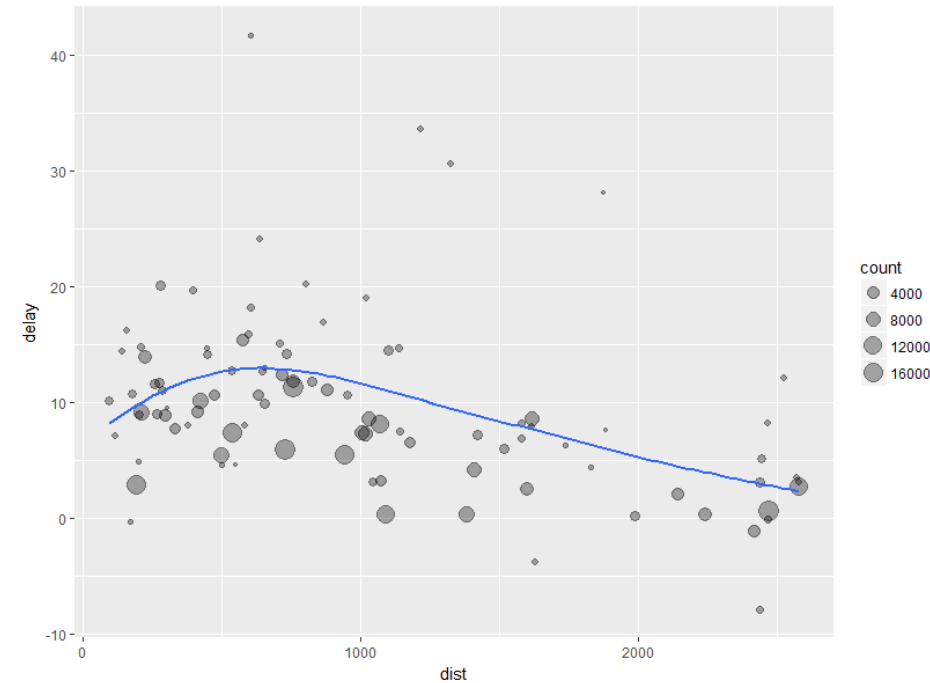
对已有列进行数据运算并添加为新列

汇总

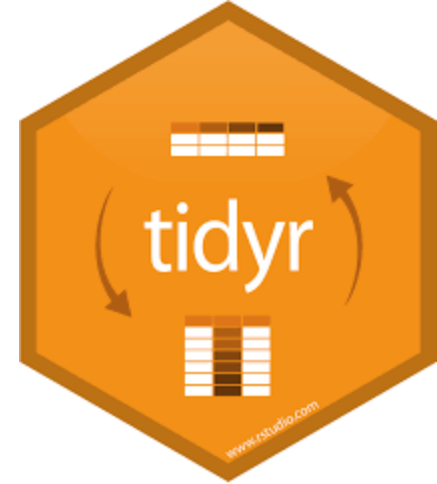
对数据框调用其它函数进行汇总操作, 返回一维的结果

管道函数(%>%) 和 绘图

```
1 delays <- flights %>%
2   group_by(dest) %>%
3   summarise(
4     count = n(),
5     dist = mean(distance, na.rm = TRUE),
6     delay = mean(arr_delay, na.rm = TRUE)
7   ) %>%
8   filter(count > 20, dest != "HNL")
9 ggplot(data = delay, mapping = aes(x = dist, y = delay)) +
10 geom_point(aes(size = count), alpha = 1/3) +
11 geom_smooth(se = FALSE)
12
```



数据整形 (Reshaping Data)



```
1 install.packages("tidyr")
2 library("tidyr")
3 # Create a new tibble use data_frame()
4 friends_data <- data_frame(
5   name = c("Nicolas", "Thierry", "Bernard", "Jerome"),
6   age = c(27, 25, 29, 26),
7   height = c(180, 170, 185, 169),
8   married = c(TRUE, FALSE, TRUE, TRUE)
9 )
10 friends_data
11 Source: local data frame [4 x 4]
12   name    age height married
13   <chr> <dbl> <dbl> <lgl>
14 1 Nicolas    27    180    TRUE
15 2 Thierry    25    170   FALSE
16 3 Bernard    29    185    TRUE
17 4 Jerome     26    169    TRUE
18 # Convert your data as a tibble
19 data("iris")
20 class(iris)
21 [1] "data.frame"
22 my_data <- as_data_frame(iris)
23 class(my_data)
24 [1] "tbl_df"      "tbl"        "data.frame"
25 # Turn a tibble back to a data frame
26 my_data2 <- as.data.frame(my_data)
27
```

Tibble Data Format in R: Best and Modern Way to Work with Data

```
> flights
```

```
Source: local data frame [336,776 x 16]
```

	year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum
	<int>	<int>	<int>	<chr>	<dbl>	<chr>	<dbl>	<chr>	<chr>
1	2013	1	1	517	2	830	11	UA	N14228
2	2013	1	1	533	4	850	20	UA	N24211
3	2013	1	1	542	2	923	33	AA	N619AA
4	2013	1	1	544	14		-18	B6	N804JB
5	2013	1	1	552	2		-25	DL	N668DN
6	2013	1	1	553	0		12	UA	N39463
7	2013	1	1	554	3		19	B6	N516JB
8	2013	1	1	555	9		-14	EV	N829AS
9	2013	1	1	556	8		-8	B6	N593JB
10	2013	1	1	558	-2	753	8	AA	N3ALAA

Column types

First 10 rows
printed by default
(Useful with large
data sets)

```
Variables not shown: flight <int>, origin <chr>, dest <chr>, air_time  
<dbl>, distance <dbl>, hour <dbl>, minute <dbl>.
```

Columns that don't
fit to the screen are
not shown

- `data_frame()`: create a new tibble data frame (tbl_df)
- `as_data_frame()`: convert your data as tbl_df



```

1 > library("tidyr")
2 > my_data <- USArrests[c(1, 10, 20, 30), ]
3 > my_data
4   |      Murder Assault UrbanPop Rape
5 Alabama      13.2    236      58 21.2
6 Georgia      17.4    211     60 25.8
7 Maryland     11.3    300     67 27.8
8 New Jersey    7.4    159     89 18.8
9 > my_data <- cbind(state = rownames(my_data), my_data)
10 > my_data
11 |      state Murder Assault UrbanPop Rape
12 Alabama  Alabama  13.2    236     58 21.2
13 Georgia  Georgia  17.4    211     60 25.8
14 Maryland Maryland 11.3    300     67 27.8
15 New Jersey New Jersey 7.4    159     89 18.8
16

```

Tidy data
Variables

State	Murder	Assault	UrbanPop	Rape	
1	State	Murder	Assault	UrbanPop	Rape
2	Alabama	13.2	236	58	21.2
3	Alaska	10	263	48	44.5
4	Arizona	8.1	294	80	31
5	Arkansas	8.8	190	50	19.5
6	California	9	276	91	40.6
7	Colorado	7.9	204	78	38.7
8	Connecticut	3.3	110	77	11.1
9	Delaware	5.9	238	72	15.8
10	Florida	15.4	335	80	31.9
11	Georgia	17.4	211	60	25.8
12	Hawaii	5.3	46	83	20.2
13	Idaho	2.6	120	54	14.2
14	Illinois	10.4	249	83	24
15	Indiana	7.2	113	65	21
16	Iowa	2.2	56	57	11.3
17	Kansas	6	115	66	18
18	Kentucky	12.9	206	43	16.3
19	Louisiana	13.2	262	41	20.3
20	Maine	5.4	111	55	11.1
21	Maryland	11.3	300	67	27.8
22	Massachusetts	8.6	182	58	15.1
23	Michigan	9.9	274	99	31.9
24	Minnesota	7.7	185	63	18.7
25	Mississippi	21.9	320	29	24.7
26	Missouri	12.6	246	76	20.1
27	Montana	1.9	67	32	10.9
28	Nebraska	6.8	131	68	16.8
29	Nevada	6.3	193	65	19.6
30	New Jersey	7.4	159	89	18.8
31	New Mexico	9.1	196	69	18.8
32	New York	9	297	109	30.1
33	North Carolina	10.8	307	108	21.6
34	North Dakota	0.6	38	37	10.7
35	Ohio	11.6	274	106	24.4
36	Oklahoma	9.5	198	63	18.5
37	Oregon	8.3	175	87	17.8
38	Rhode Island	7.1	114	52	15.1
39	Texas	17.1	307	119	30.1
40	Vermont	5.1	85	42	11.1
41	Virginia	10.3	182	72	18.3
42	Washington	7.6	158	75	17.6
43	West Virginia	17.4	307	29	24.7
44	Wisconsin	9.7	187	78	19.7
45	Wyoming	3.1	69	33	10.1

tidy data set:

- each column represents a variable
- each row represents an observation
- ❖ The opposite of tidy is messy data

Organize Your Data for Easier Analyses in R

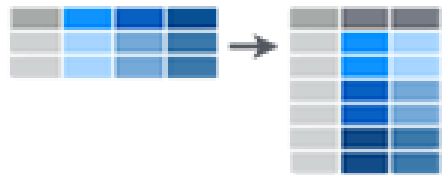


- **gather():** collapse multiple columns into key-value pairs
- **spread():** reverse of gather. Separate one column into multiple
- **separate():** separate one column into multiple
- **unite():** unite multiple columns into one

tidyr R package

- ✓ gather和spread函数将数据在长格式和宽格式之间相互转化，应用在比如稀疏矩阵和稠密矩阵之间的转化；
- ✓ separate和union方法提供了数据分组拆分、合并的功能，应用在nominal数据的转化上；

gather() 宽数据转为长数据



`gather(data, key, value, ...)`

- **data:** A data frame
- **key, value:** Names of key and value columns to create in output
- **...:** Specification of columns to gather. Allowed values are:
 - variable names
 - if you want to select all variables between a and e, use a:e
 - if you want to exclude a column name y use -y
 - for more options, see: `dplyr::select()`

```
1 my_data2 <- gather(my_data,
2   key = "arrest_attribute",
3   value = "arrest_estimate",
4   -state)
5 # Only Murder and Assault columns,
6 # remaining columns have been duplicated.
7 my_data2 <- gather(my_data,
8   key = "arrest_attribute",
9   value = "arrest_estimate",
10  Murder, Assault)
11
```

```
1 > my_data
2   state Murder Assault UrbanPop Rape
3 Alabama  13.2   236     58 21.2
4 Georgia  17.4   211     60 25.8
5 Maryland 11.3   300     67 27.8
6 New Jersey 7.4   159     89 18.8
7 -----
8 > my_data2
9   state arrest_attribute arrest_estimate
10 1 Alabama Murder 13.2
11 2 Georgia Murder 17.4
12 3 Maryland Murder 11.3
13 4 New Jersey Murder 7.4
14 5 Alabama Assault 236.0
15 6 Georgia Assault 211.0
16 7 Maryland Assault 300.0
17 8 New Jersey Assault 159.0
18 9 Alabama UrbanPop 58.0
19 10 Georgia UrbanPop 60.0
20 11 Maryland UrbanPop 67.0
21 12 New Jersey UrbanPop 89.0
22 13 Alabama Rape 21.2
23 14 Georgia Rape 25.8
24 15 Maryland Rape 27.8
25 16 New Jersey Rape 18.8
26
```

spread() 长数据转为宽数据



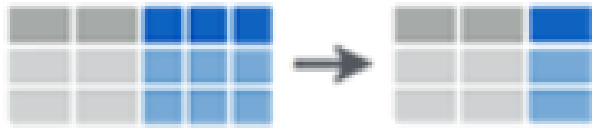
spread(data, key, value)

- **data:** A data frame
- **key:** The (unquoted) name of the column whose values will be used as column headings.
- **value:** The (unquoted) names of the column whose values will populate the cells.

```
1 my_data3 <- spread(my_data2,  
2   key = "arrest_attribute",  
3   value = "arrest_estimate"  
4 )  
5
```

```
1 > my_data  
2   state Murder Assault UrbanPop Rape  
3 Alabama Alabama 13.2 236 58 21.2  
4 Georgia Georgia 17.4 211 60 25.8  
5 Maryland Maryland 11.3 300 67 27.8  
6 New Jersey New Jersey 7.4 159 89 18.8  
7 -----  
8 > my_data2  
9   state arrest_attribute arrest_estimate  
10 1 Alabama Murder 13.2  
11 2 Georgia Murder 17.4  
12 3 Maryland Murder 11.3  
13 4 New Jersey Murder 7.4  
14 5 Alabama Assault 236.0  
15 6 Georgia Assault 211.0  
16 7 Maryland Assault 300.0  
17 8 New Jersey Assault 159.0  
18 9 Alabama UrbanPop 58.0  
19 10 Georgia UrbanPop 60.0  
20 11 Maryland UrbanPop 67.0  
21 12 New Jersey UrbanPop 89.0  
22 13 Alabama Rape 21.2  
23 14 Georgia Rape 25.8  
24 15 Maryland Rape 27.8  
25 16 New Jersey Rape 18.8  
26 -----  
27 > my_data3  
28   state Assault Murder Rape UrbanPop  
29 1 Alabama 236 13.2 21.2 58  
30 2 Georgia 211 17.4 25.8 60  
31 3 Maryland 300 11.3 27.8 67  
32 4 New Jersey 159 7.4 18.8 89  
33
```


unite() 多列合并为一列



```
unite(data, col, ..., sep = "_")
```

- ✓ **data**: A data frame
- ✓ **col**: The new (unquoted) name of column to add.
- ✓ **sep**: Separator to use between values

```
1 my_data4 <- unite(my_data,  
2   col = "Murder_Assault",  
3   Murder, Assault,  
4   sep = "_")  
5
```

```
1 > my_data  
2   Murder Assault UrbanPop Rape  
3 Alabama    13.2    236     58 21.2  
4 Georgia    17.4    211     60 25.8  
5 Maryland   11.3    300     67 27.8  
6 New Jersey  7.4    159     89 18.8  
7 > my_data4  
8   Murder_Assault UrbanPop Rape  
9 Alabama    13.2_236     58 21.2  
10 Georgia    17.4_211     60 25.8  
11 Maryland   11.3_300     67 27.8  
12 New Jersey  7.4_159     89 18.8  
13
```

separate() 将一列分离为多列



```
separate(data, col, into, sep = "[^[:alnum:]]+")
```

- ✓ **data:** A data frame
- ✓ **col:** Unquoted column names
- ✓ **into:** Character vector specifying the names of new variables to be created.
- ✓ **sep:** Separator between columns:
 - ✓ If character, is interpreted as a regular expression.
 - ✓ If numeric, interpreted as positions to split at. Positive values start at 1 at the far-left of the string; negative value start at -1 at the far-right of the string.

```
1 my_data5 <- separate(my_data4,  
2   col = "Murder_Assault",  
3   into = c("Murder", "Assault"),  
4   sep = "_")  
5
```

```
1 > my_data  
2   Murder Assault UrbanPop Rape  
3 Alabama      13.2    236     58 21.2  
4 Georgia      17.4    211     60 25.8  
5 Maryland     11.3    300     67 27.8  
6 New Jersey   7.4     159     89 18.8  
7 > my_data4  
8   Murder_Assault UrbanPop Rape  
9 Alabama      13.2_236     58 21.2  
10 Georgia     17.4_211     60 25.8  
11 Maryland    11.3_300     67 27.8  
12 New Jersey  7.4_159     89 18.8  
13 > my_data5  
14  Murder Assault UrbanPop Rape  
15 Alabama      13.2    236     58 21.2  
16 Georgia      17.4    211     60 25.8  
17 Maryland     11.3    300     67 27.8  
18 New Jersey   7.4     159     89 18.8  
19
```

管道函数(%>%)

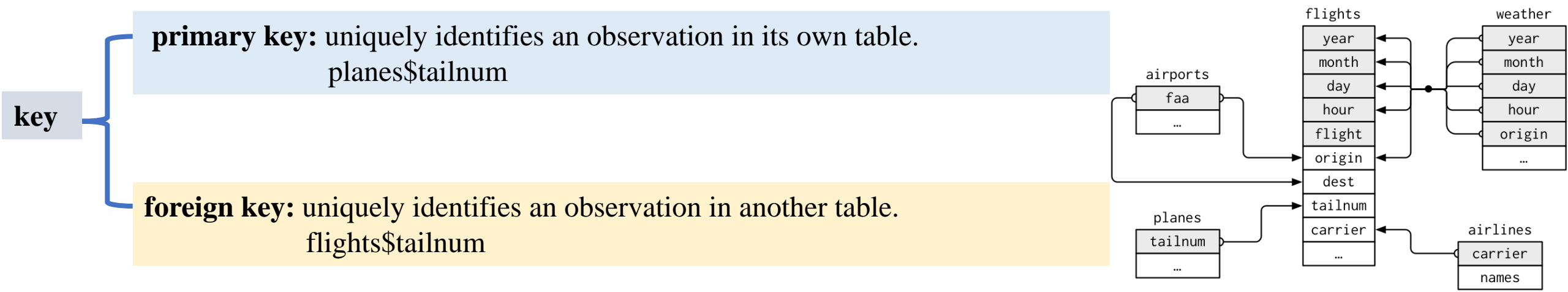
gather() + unite()

```
1 my_data6 <- my_data %>% gather(key = "arrest_attribute",
2 |                               value = "arrest_estimate",
3 |                               Murder:UrbanPop) %>%
4 unite(col = "attribute_estimate",
5 |      arrest_attribute, arrest_estimate)
6 -----
7 > my_data
8 |           Murder Assault UrbanPop Rape
9 | Alabama      13.2    236      58 21.2
10 | Georgia      17.4    211      60 25.8
11 | Maryland     11.3    300      67 27.8
12 | New Jersey   7.4     159      89 18.8
13 > my_data6
14 | Rape attribute_estimate
15 | 1 21.2      Murder_13.2
16 | 2 25.8      Murder_17.4
17 | 3 27.8      Murder_11.3
18 | 4 18.8      Murder_7.4
19 | 5 21.2      Assault_236
20 | 6 25.8      Assault_211
21 | 7 27.8      Assault_300
22 | 8 18.8      Assault_159
23 | 9 21.2      UrbanPop_58
24 | 10 25.8     UrbanPop_60
25 | 11 27.8     UrbanPop_67
26 | 12 18.8     UrbanPop_89
27
```

关系型数据 (Relational data)



```
1 library(nycflights13)
2 > class(flights)
3 [1] "tbl_df"      "tbl"        "data.frame"
4 > flights
5 # A tibble: 336,776 x 19
6   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier flight tailnum origin dest
7   <int> <int> <int> <int>         <int>         <dbl>   <int>         <int>         <dbl>   <chr>   <int>   <chr>   <chr> <chr>
8     1  2013     1     1     517             515           2         830             819          11     UA    1545   N14228   EWR    IAH
9     2  2013     1     1     533             529           4         850             830          20     UA    1714   N24211   LGA    IAH
10    3  2013     1     1     542             540           2         923             850          33     AA    1141   N619AA   JFK    MIA
11    4  2013     1     1     544             545          -1        1004            1022         -18     B6     725   N804JB   JFK    BQN
12    5  2013     1     1     554             600          -6         812             837         -25     DL     461   N668DN   LGA    ATL
13    6  2013     1     1     554             558          -4         740             728          12     UA    1696   N39463   EWR    ORD
14    7  2013     1     1     555             600          -5         913             854          19     B6     507   N516JB   EWR    FLL
15    8  2013     1     1     557             600          -3         709             723         -14     EV    5708   N829AS   LGA    IAD
16    9  2013     1     1     557             600          -3         838             846          -8     B6      79   N593JB   JFK    MCO
17   10  2013     1     1     558             600          -2         753             745           8     AA     301   N3ALAA   LGA    ORD
18 # ... with 336,766 more rows, and 5 more variables: air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
19 planes %>% count(tailnum) %>% filter(n > 1) # verify the key
20
```

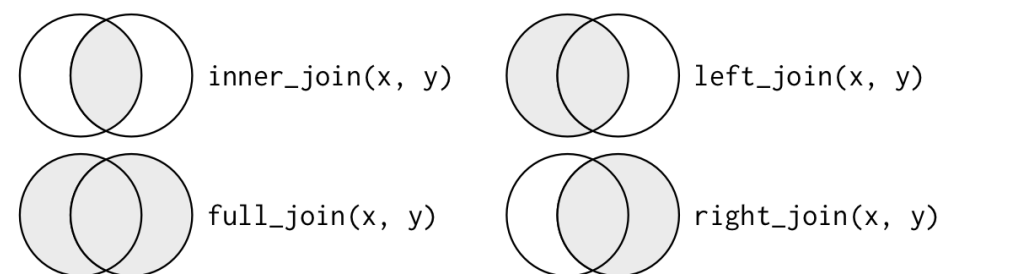
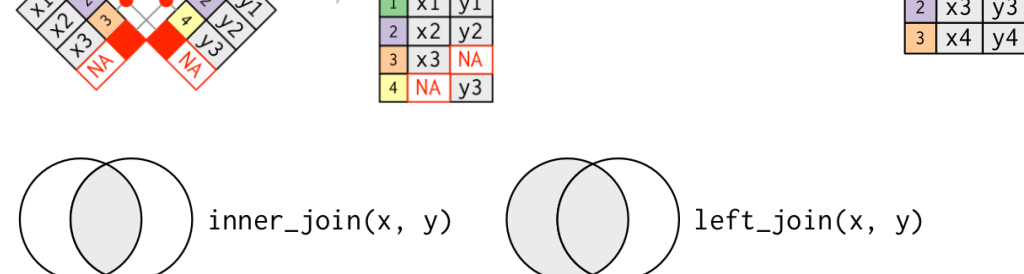
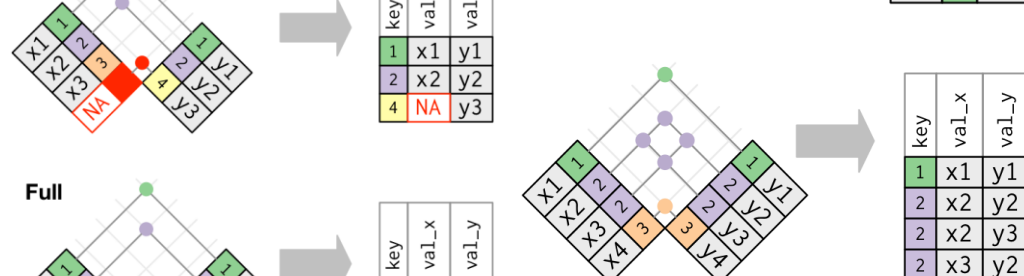
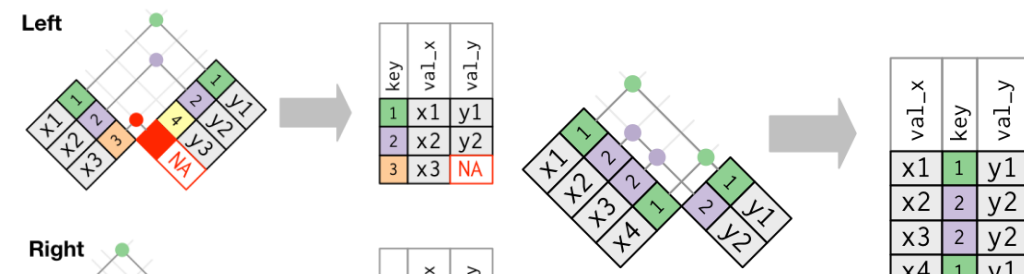
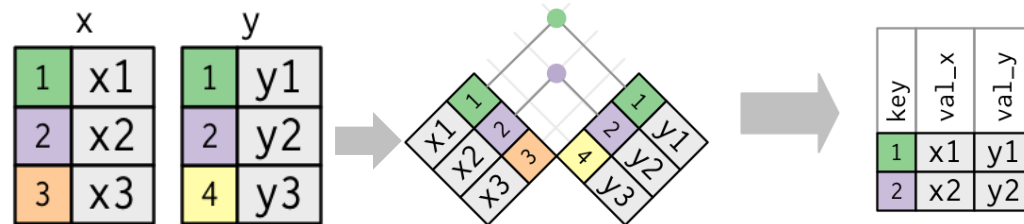


Mutating joins

```

1 > flights2 <- flights %>% select(year:day, hour, origin, dest, tailnum, carrier)
2 > flights2
3 # A tibble: 336,776 x 8
4   year month   day hour origin dest tailnum carrier
5   <int> <int> <int> <dbl> <chr> <chr> <chr> <chr>
6 1 2013     1     1     5 EWR  IAH  N14228  UA
7 2 2013     1     1     5 LGA  IAH  N24211  UA
8 3 2013     1     1     5 JFK  MIA  N619AA  AA
9 4 2013     1     1     5 JFK  BQN  N804JB  B6
10 5 2013     1     1     6 LGA  ATL  N668DN  DL
11 6 2013     1     1     5 EWR  ORD  N39463  UA
12 7 2013     1     1     6 EWR  FLL  N516JB  B6
13 8 2013     1     1     6 LGA  IAD  N829AS  EV
14 9 2013     1     1     6 JFK  MCO  N593JB  B6
15 10 2013     1     1     6 LGA  ORD  N3ALAA  AA
16 # ... with 336,766 more rows
17 > flights2 %>% select(-origin, -dest) %>% left_join(airlines, by = "carrier")
18 # A tibble: 336,776 x 7
19   year month   day hour tailnum carrier      name
20   <int> <int> <int> <dbl> <chr> <chr> <chr>
21 1 2013     1     1     5 N14228  UA      United Air Lines Inc.
22 2 2013     1     1     5 N24211  UA      United Air Lines Inc.
23 3 2013     1     1     5 N619AA  AA      American Airlines Inc.
24 4 2013     1     1     5 N804JB  B6      JetBlue Airways
25 5 2013     1     1     6 N668DN  DL      Delta Air Lines Inc.
26 6 2013     1     1     5 N39463  UA      United Air Lines Inc.
27 7 2013     1     1     6 N516JB  B6      JetBlue Airways
28 8 2013     1     1     6 N829AS  EV      ExpressJet Airlines Inc.
29 9 2013     1     1     6 N593JB  B6      JetBlue Airways
30 10 2013     1     1     6 N3ALAA  AA      American Airlines Inc.
31 # ... with 336,766 more rows
32

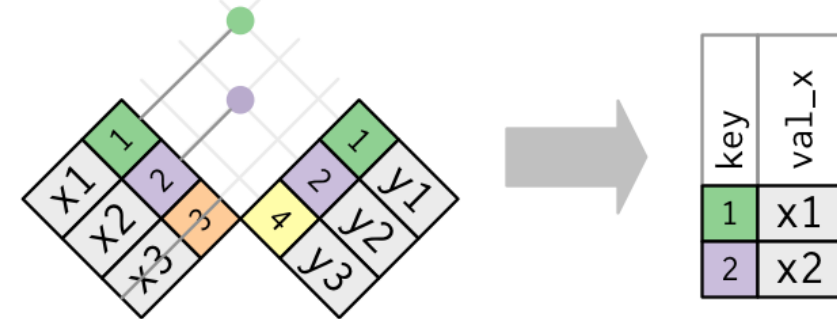
```



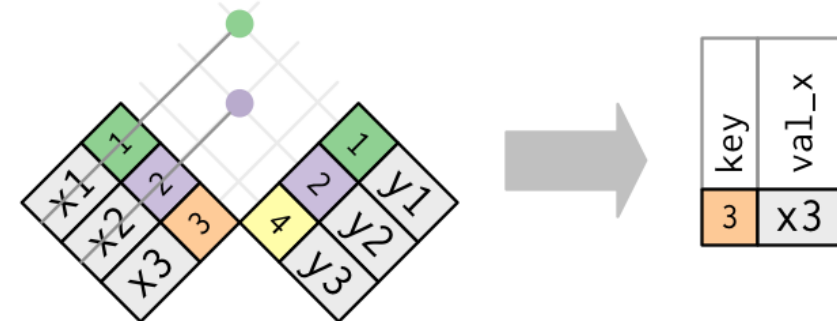
Filtering joins

```
1 > top_dest <- flights %>% count(dest, sort = TRUE) %>% head(10)
2 > flights %>% semi_join(top_dest)
3 > dim(flights)
4 [1] 336776    19
5 > dim(flights %>% semi_join(top_dest))
6 [1] 80262    19
7 > flights %>% anti_join(planes, by = "tailnum") %>% count(tailnum, sort = TRUE)
8 # A tibble: 722 x 2
9   tailnum     n
10  <chr> <int>
11 1 <NA> 2512
12 2 N725MQ 575
13 3 N722MQ 513
14 4 N723MQ 507
15 5 N713MQ 483
16 6 N735MQ 396
17 # ... with 712 more rows
18
```

semi_join(x, y): keeps all observations in x that have a match in y



anti_join(x, y): drops all observations in x that have a match in y



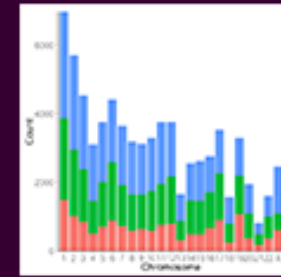
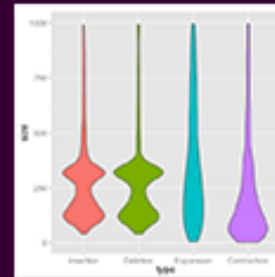
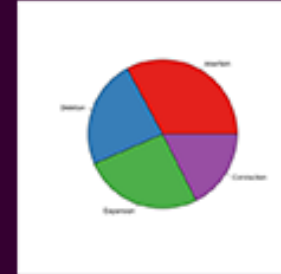
Set operations

- ✓ **intersect(x, y)**: return only observations in both x and y.
- ✓ **union(x, y)**: return unique observations in x and y.
- ✓ **setdiff(x, y)**: return observations in x, but not in y.

```
1 df1 <- tribble(
2   ~x, ~y,
3   1,  1,
4   2,  1
5 )
6 df2 <- tribble(
7   ~x, ~y,
8   1,  1,
9   1,  2
10 )
11 > df1
12 # A tibble: 2 x 2
13   x     y
14 <dbl> <dbl>
15 1     1     1
16 2     2     1
17 > df2
18 # A tibble: 2 x 2
19   x     y
20 <dbl> <dbl>
21 1     1     1
22 2     1     2
23
```

```
1 > intersect(df1, df2)
2 # A tibble: 1 x 2
3   x     y
4 <dbl> <dbl>
5 1     1     1
6 -----
7 > union(df1, df2)
8 # A tibble: 3 x 2
9   x     y
10 <dbl> <dbl>
11 1     1     2
12 2     2     1
13 3     1     1
14 -----
15 > setdiff(df1, df2)
16 # A tibble: 1 x 2
17   x     y
18 <dbl> <dbl>
19 1     2     1
20
```

Plotting in R for Biologists



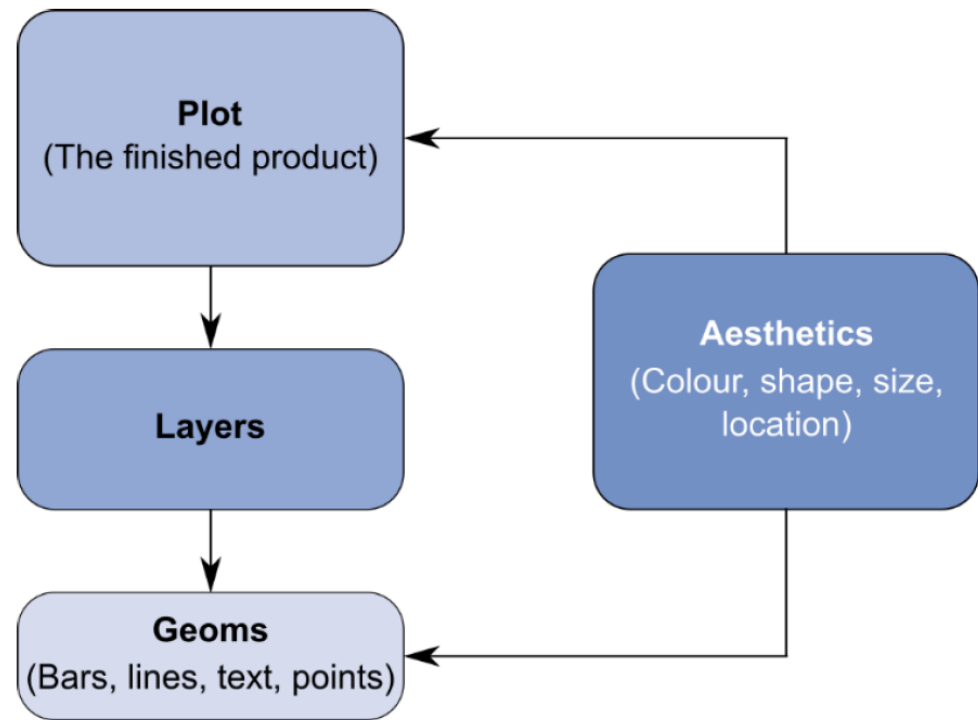
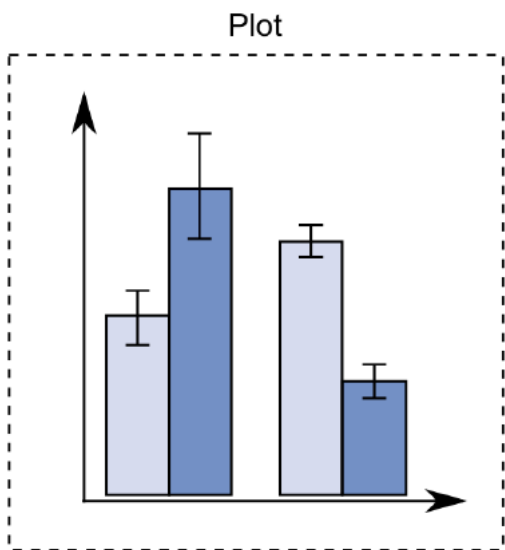
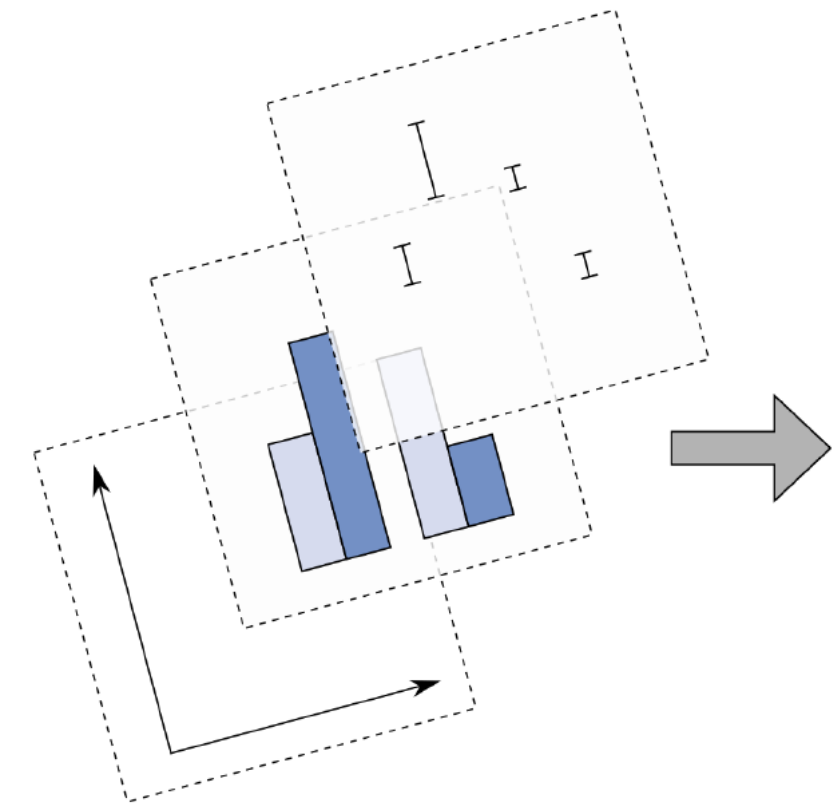
ggplot2

<http://hadley.nz/>

<http://ggplot2.org/>



图层 (layers)



ggplot2的基本概念

➤ 数据 (Data) 和映射 (Mapping)

length	width	depth	trt
2	3	4	a
1	2	1	a
4	5	15	b
9	10	80	b



Mapping

x	y	colour
2	3	a
1	2	a
4	5	b
9	10	b



Scale

x	y	colour
25	11	red
0	0	red
75	53	blue
200	300	blue

➤ 标度 (Scale)

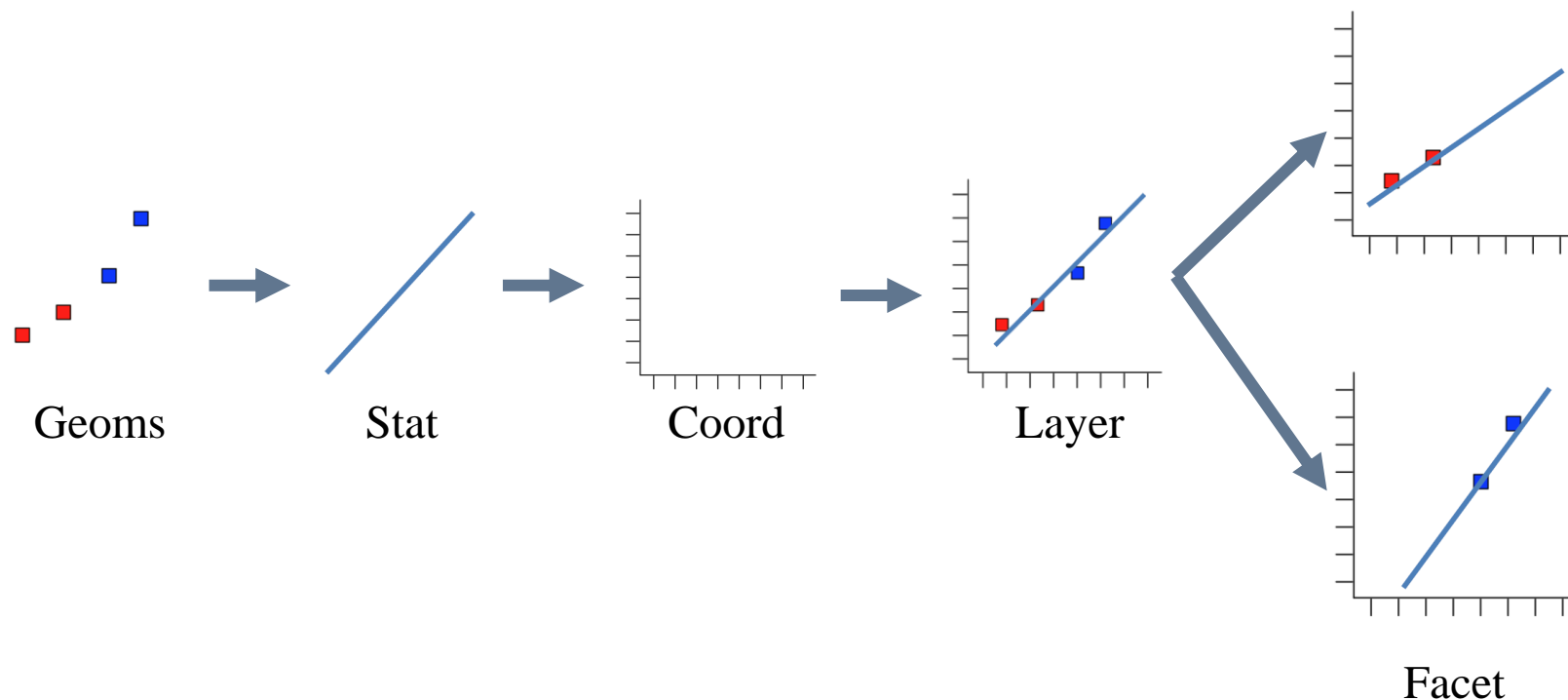
➤ 几何对象 (Geometric)

➤ 统计变换 (Statistics)

➤ 坐标系 (Coordinate)

➤ 图层 (Layer)

➤ 分面 (Facet)



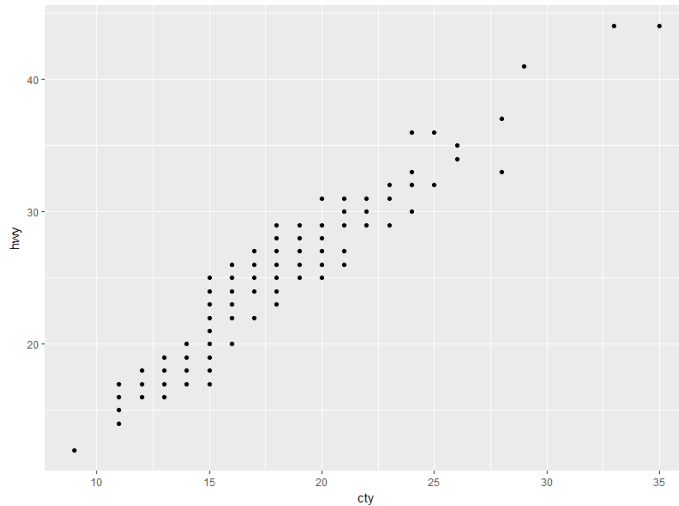
实例数据

```
1 > library(ggplot2)
2 > str(mpg)
3 	Classes 'tbl_df', 'tbl' and 'data.frame':	 234 obs. of  11 variables:
4 	 $ manufacturer: chr  "audi" "audi" "audi" "audi" ...
5 	 $ model       : chr  "a4" "a4" "a4" "a4" ...
6 	 $ displ      : num  1.8 1.8 2 2 2.8 2.8 3.1 1.8 1.8 2 ...
7 	 $ year       : int  1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 ...
8 	 $ cyl        : int  4 4 4 4 6 6 6 4 4 4 ...
9 	 $ trans      : chr  "auto(l5)" "manual(m5)" "manual(m6)" "auto(av)" ...
10 	 $ drv       : chr  "f" "f" "f" "f" ...
11 	 $ cty       : int  18 21 20 21 16 18 18 18 16 20 ...
12 	 $ hwy      : int  29 29 31 30 26 26 27 26 25 28 ...
13 	 $ fl       : chr  "p" "p" "p" "p" ...
14 	 $ class    : chr  "compact" "compact" "compact" "compact" ...
15
16 > head(mpg)
17 	# A tibble: 6 x 11
18 	 manufacturer model displ year cyl trans drv cty hwy fl class
19 	 <chr> <chr> <dbl> <int> <int> <chr> <chr> <int> <int> <chr> <chr>
20 	 1 audi a4 1.8 1999 4 auto(l5) f 18 29 p compact
21 	 2 audi a4 1.8 1999 4 manual(m5) f 21 29 p compact
22 	 3 audi a4 2.0 2008 4 manual(m6) f 20 31 p compact
23 	 4 audi a4 2.0 2008 4 auto(av) f 21 30 p compact
24 	 5 audi a4 2.8 1999 6 auto(l5) f 16 26 p compact
25 	 6 audi a4 2.8 1999 6 manual(m5) f 18 26 p compact
26
```


散点图

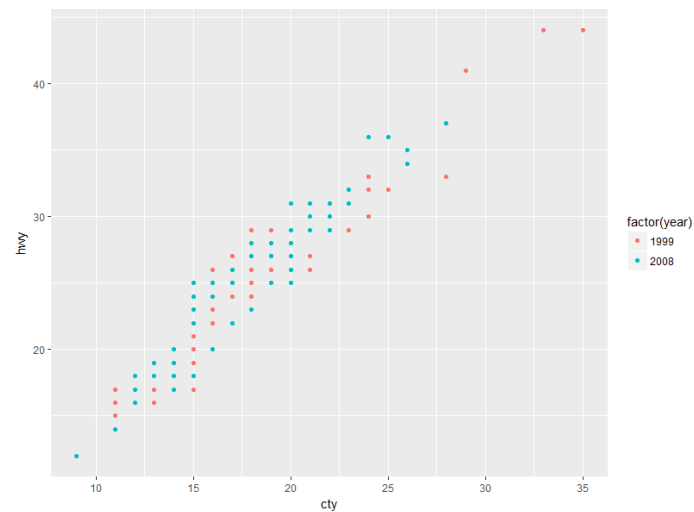
aesthetics

```
1 p <- ggplot(data=mpg, mapping=aes(x=cty, y=hwy))  
2 p + geom_point()  
3
```



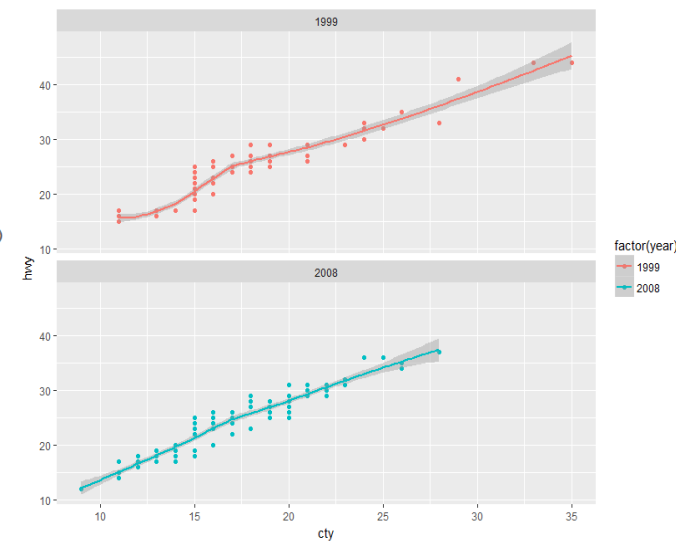
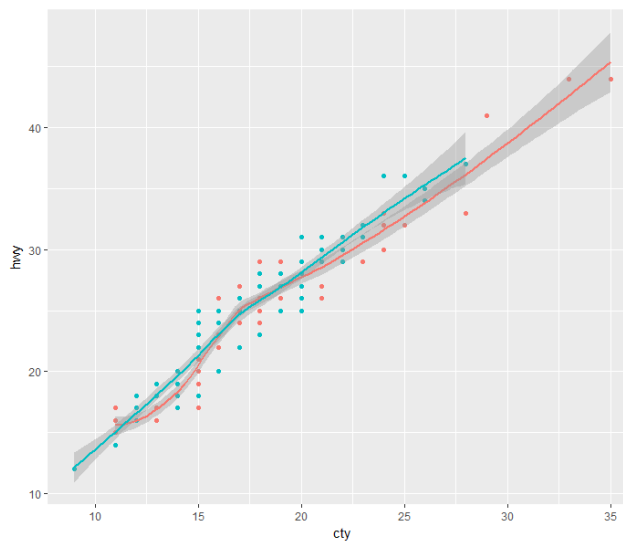
将年份映射到颜色属性 (Mapping):

```
1 p <- ggplot(mpg, aes(x=cty, y=hwy, colour=factor(year)))  
2 p + geom_point()  
3
```



增加平滑曲线 (Statistics):

```
1 p + geom_point() + stat_smooth()  
2 > summary(p + geom_point() + stat_smooth())  
3 geom_smooth: na.rm = FALSE  
4 stat_smooth: method = auto, formula = y ~ x,  
5 se = TRUE, n = 80, fullrange = FALSE,  
6 level = 0.95, na.rm = FALSE,  
7 method.args = list(), span = 0.75  
8 position_identity  
9
```

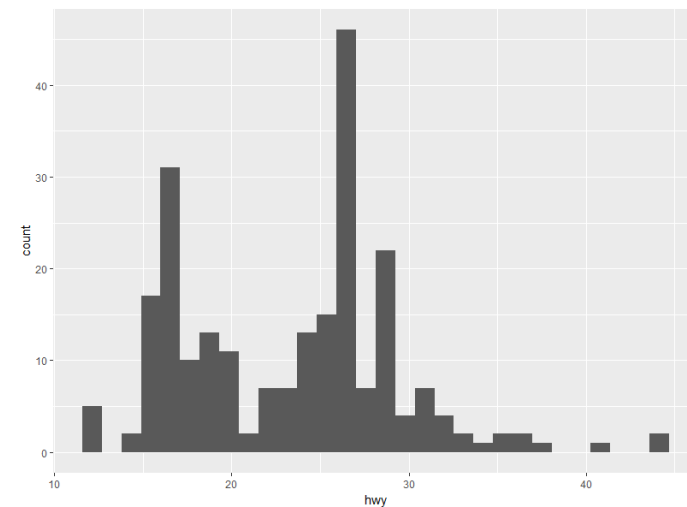


分面 (Facet):

```
1 p + geom_point() + stat_smooth()+facet_wrap(~ year, ncol=1)  
2
```

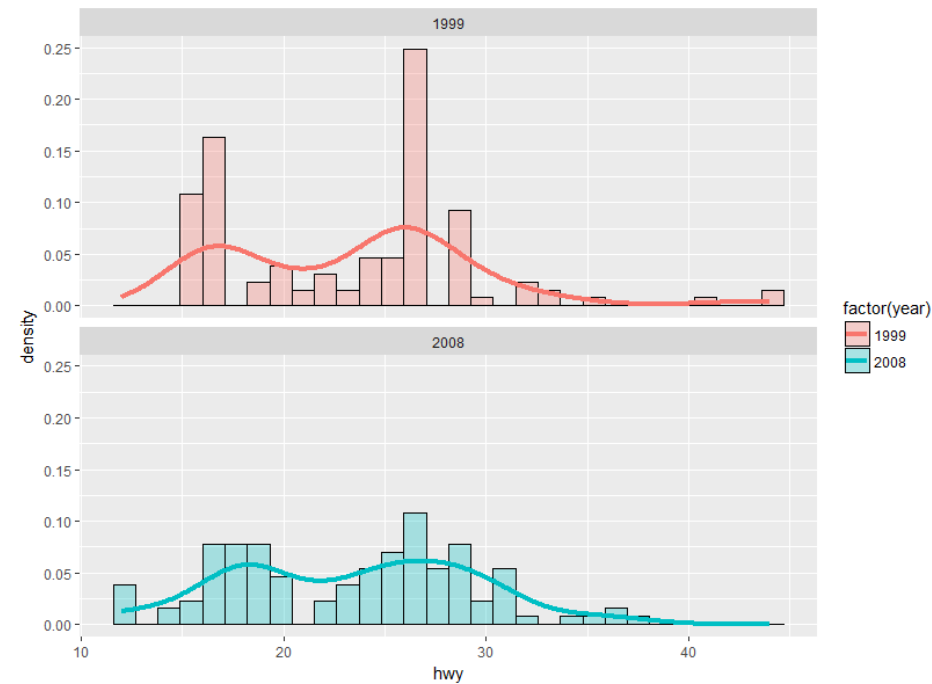
直方图

```
1 p <- ggplot(mpg, aes(x=hwy))  
2 p + geom_histogram()  
3
```



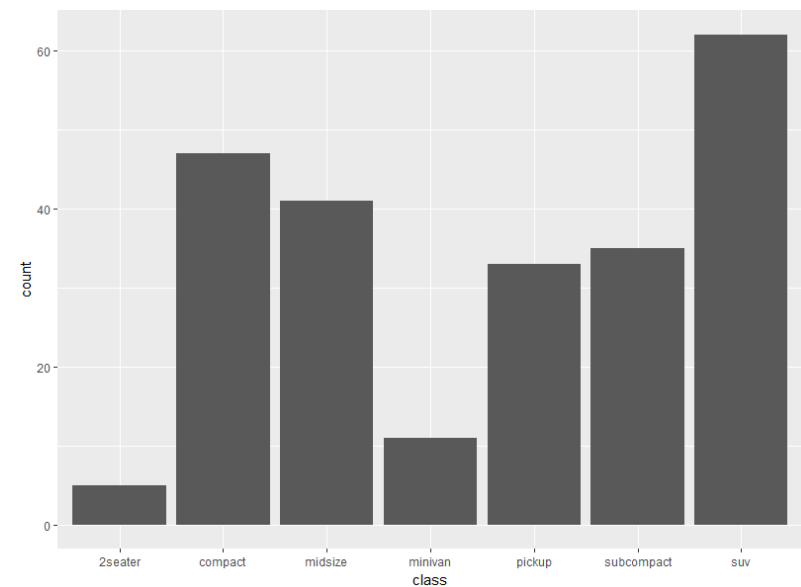
统计变换 + 分面:

```
1 # Histogram with density plot  
2 p + geom_histogram(aes(fill=factor(year), y=..density..), alpha=0.3, colour='black') +  
3   stat_density(geom='line', position='identity', size=1.5, aes(colour=factor(year))) +  
4   facet_wrap(~year, ncol=1)  
5
```



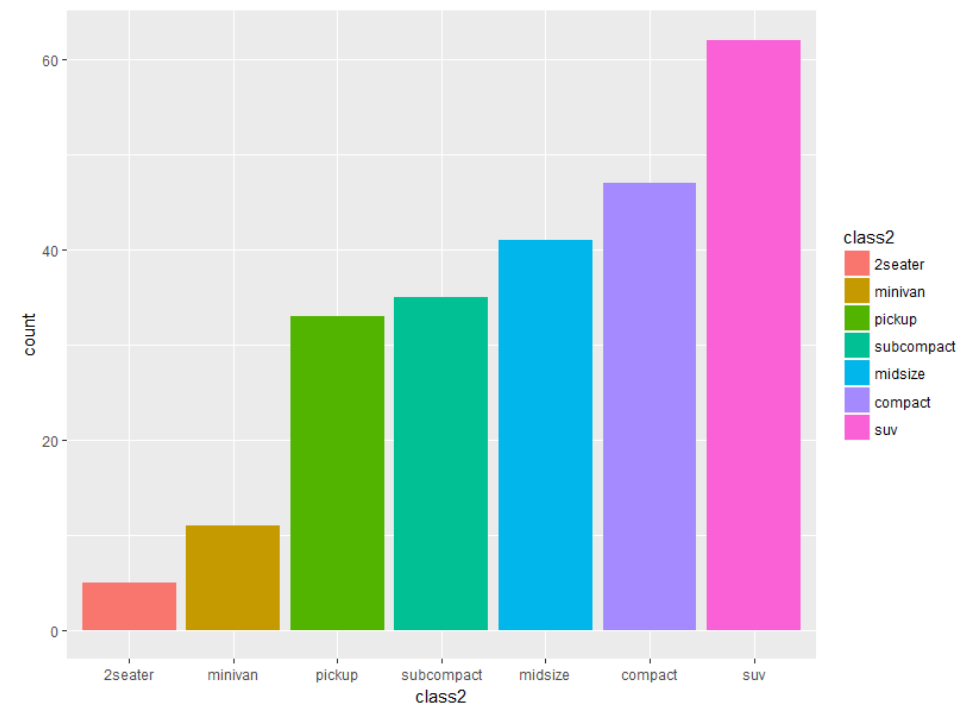
条形图

```
1 p <- ggplot(mpg, aes(x=class))
2 p + geom_bar()
3
```



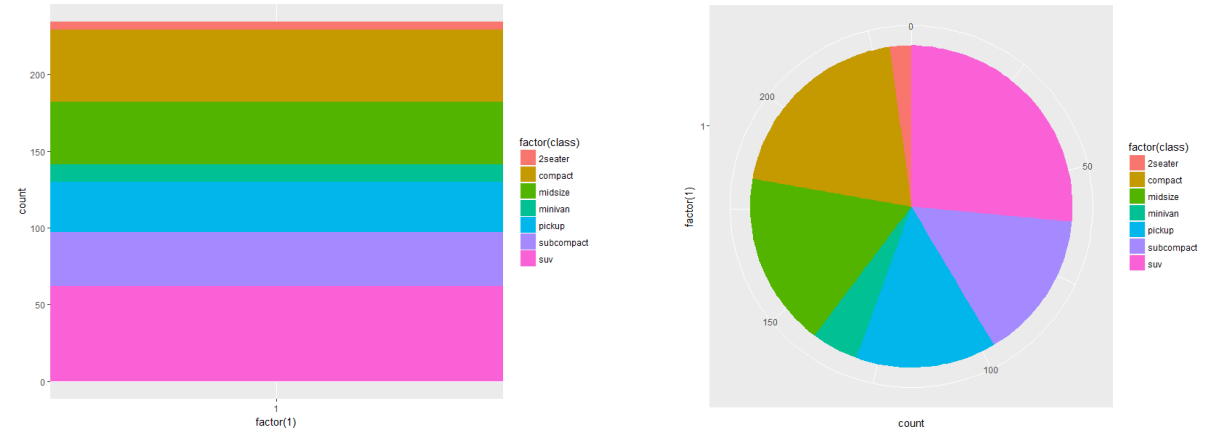
根据计数排序后绘制的条形图:

```
1 class2 <- mpg$class
2 class2 <- reorder(class2,class2,length)
3 mpg$class2 <- class2
4 p <- ggplot(mpg, aes(x=class2))
5 p + geom_bar(aes(fill=class2))
6 -----
7 help(reorder)
8
9 Reorder Levels of a Factor
10 Usage
11 reorder(x, X, FUN = mean, ...,
12         order = is.ordered(x))
13 -----
14 class(class2)
15 [1] "factor"
16
```



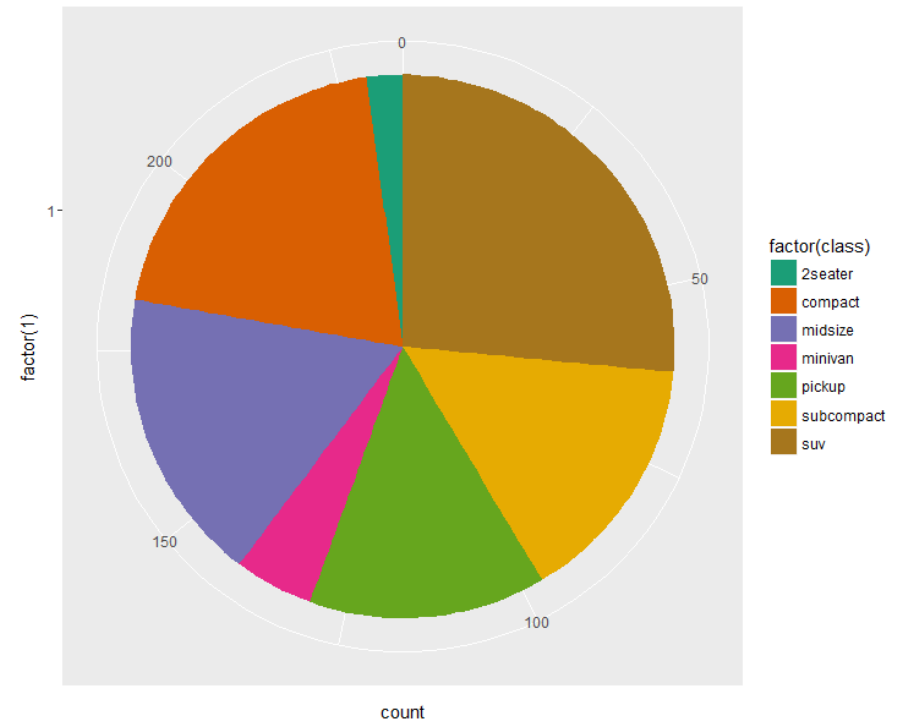
饼图

```
1 p <- ggplot(mpg, aes(x = factor(1), fill = factor(class))) +  
2   geom_bar(width = 1)  
3 p  
4  
5  
6 p + coord_polar(theta = "y")  
7
```



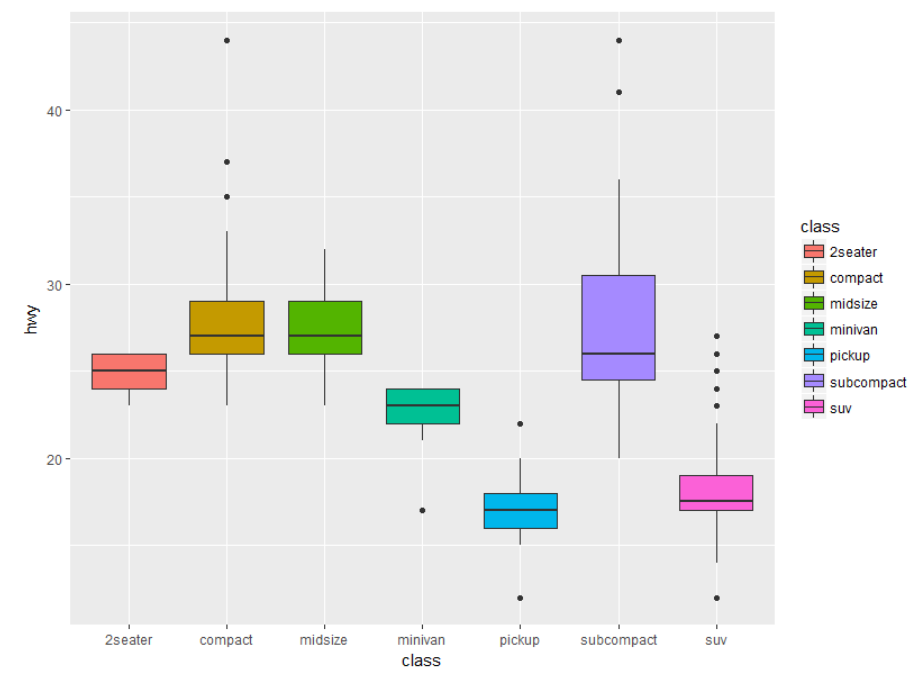
改变填充颜色:

```
1 p + coord_polar(theta = "y") +  
2   scale_fill_brewer(palette="Dark2")  
3
```



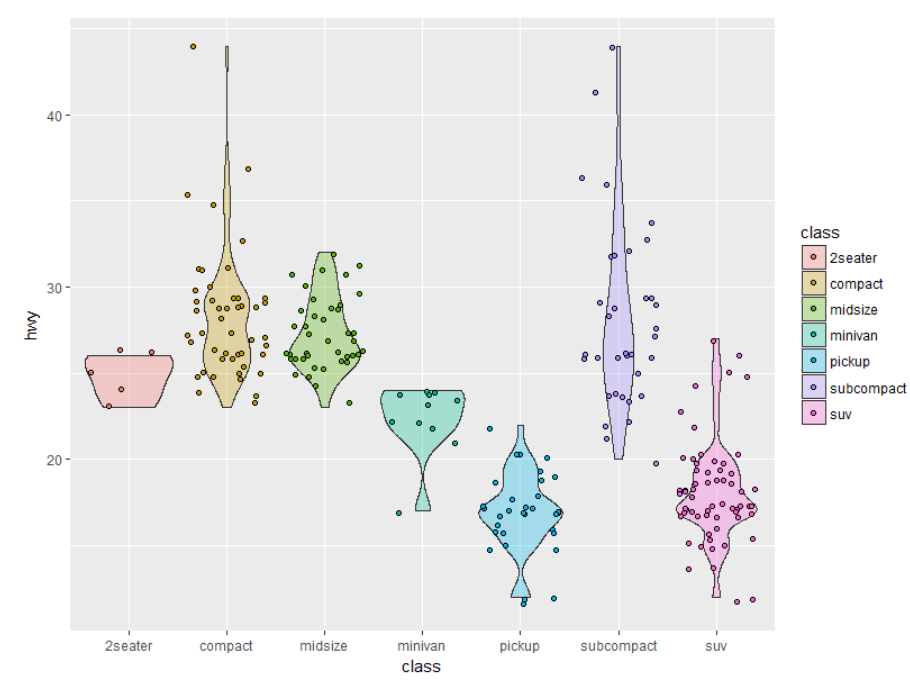
箱线图

```
1 p <- ggplot(mpg, aes(class, hwy, fill=class))  
2 p + geom_boxplot()  
3
```



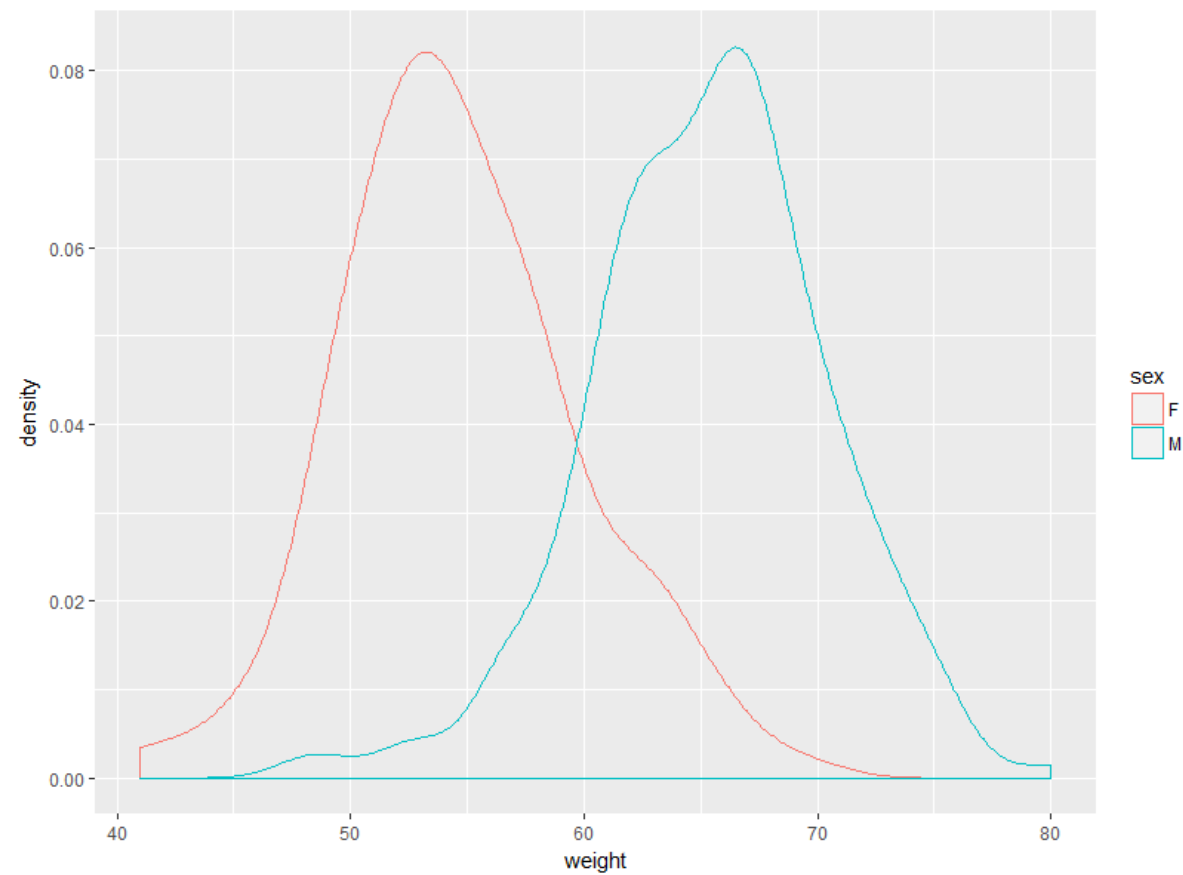
小提琴图

```
1 p + geom_violin(alpha=0.3, width=0.9) +  
2   geom_jitter(shape=21)  
3
```



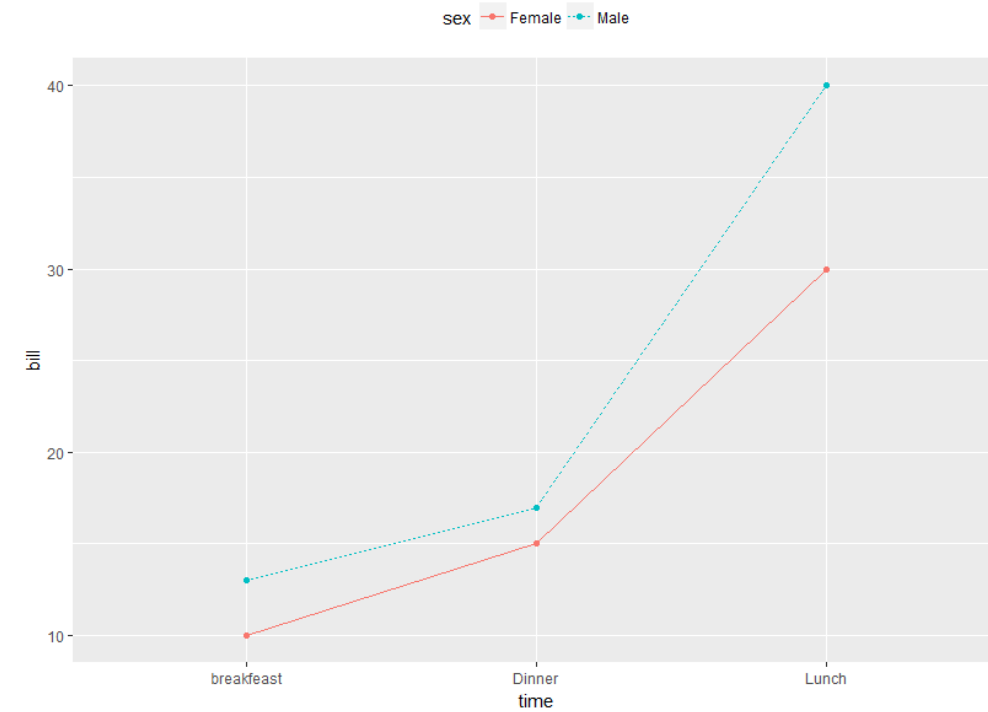
密度图

```
1 set.seed(1234)
2 df <- data.frame(
3   sex=factor(rep(c("F", "M"), each=200)),
4   weight=round(c(rnorm(200, mean=55, sd=5),
5                 rnorm(200, mean=65, sd=5)))
6 )
7 head(df)
8   sex weight
9   1    F    49
10  2    F    56
11  3    F    60
12  4    F    43
13  5    F    57
14  6    F    58
15 p <- ggplot(df, aes(x=weight, color=sex)) +
16   geom_density()
17 p
18
```



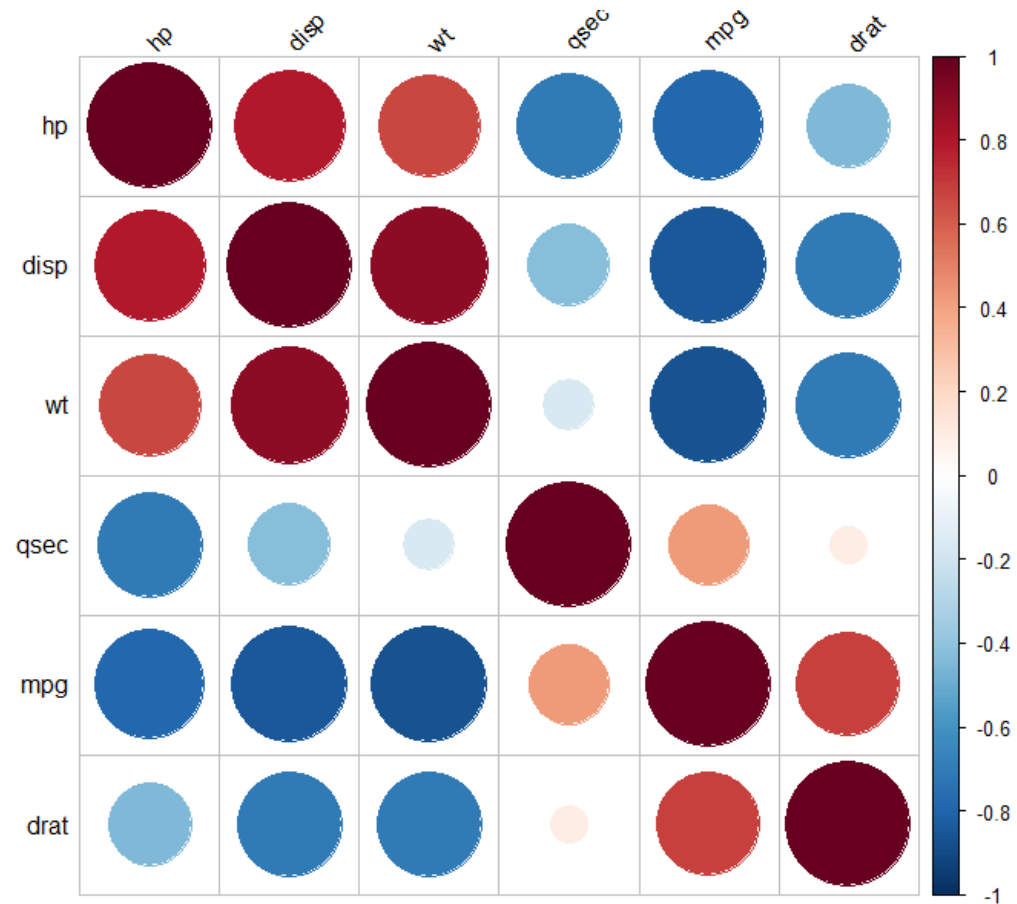
线图

```
1 df2 <- data.frame(sex = rep(c("Female", "Male"), each=3),
2                   time=c("breakfast", "Lunch", "Dinner"),
3                   bill=c(10, 30, 15, 13, 40, 17) )
4 head(df2)
5   sex      time bill
6 1 Female breakfast  10
7 2 Female      Lunch  30
8 3 Female      Dinner  15
9 4  Male breakfast  13
10 5  Male      Lunch  40
11 6  Male      Dinner  17
12 P <- ggplot(df2, aes(x=time, y=bill, group=sex)) +
13   geom_line(aes(linetype=sex, color=sex))+
14   geom_point(aes(color=sex))+
15   theme(legend.position="top")
16 P
17
```



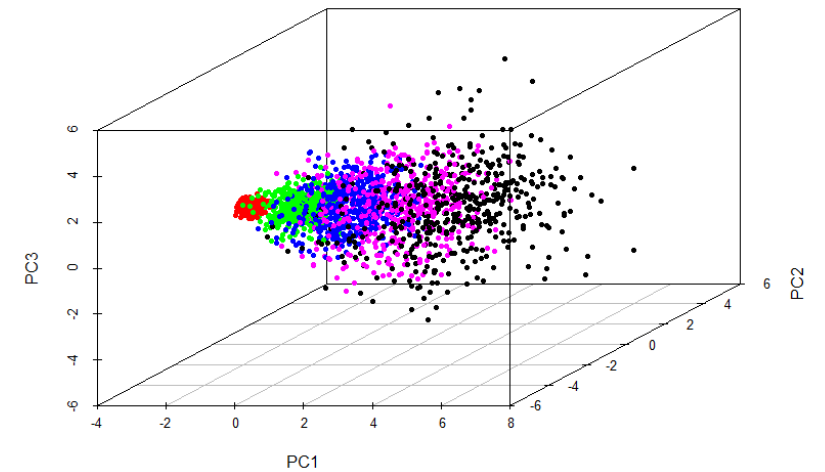
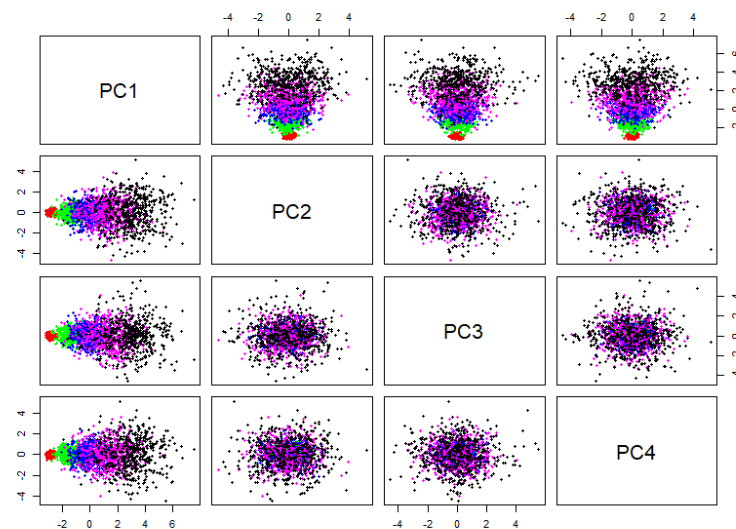
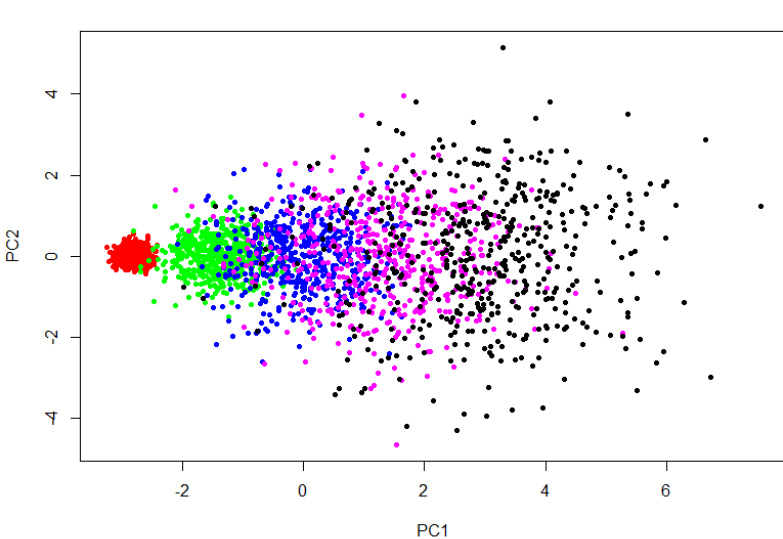
相关性图

```
1 install.packages("corrplot")
2 library(corrplot)
3 mydata <- select(mtcars, hp, disp, wt, qsec, mpg, drat)
4 source("http://www.sthda.com/upload/rquery_cormat.r")
5 rquery.cormat<-function(x, type=c('lower', 'upper', 'full', 'flatten'),...
8 { ...
91 }
92 rquery.cormat(mydata, type="full")
93
```



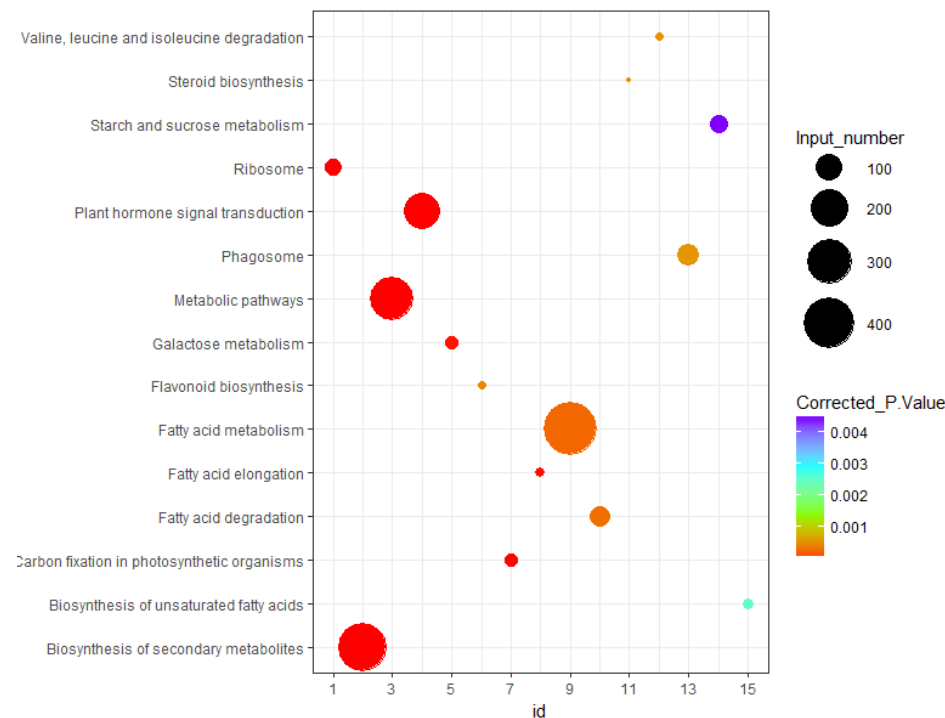
主成份分析 (PCA)

```
1  z1 <- rnorm(10000, mean=1, sd=1)
2  z2 <- rnorm(10000, mean=3, sd=3)
3  z3 <- rnorm(10000, mean=5, sd=5)
4  z4 <- rnorm(10000, mean=7, sd=7)
5  z5 <- rnorm(10000, mean=9, sd=9)
6  mydata <- matrix(c(z1, z2, z3, z4, z5), 2500, 20, byrow=T, dimnames=list(paste("R", 1:2500, sep=""), paste("C", 1:20, sep="")))
7  pca <- prcomp(mydata, scale=T)
8  summary(pca)$importance[, 1:6]
9  mycolors <- c("red", "green", "blue", "magenta", "black")
10 -----
11 plot(pca$x[,1:2], pch=20, col=mycolors[sort(rep(1:5, 500))])
12 -----
13 pairs(pca$x[,1:4], pch=20, col=mycolors[sort(rep(1:5, 500))])
14 -----
15 library(scatterplot3d)
16 scatterplot3d(pca$x[,1:3], pch=20, color=mycolors[sort(rep(1:5, 500))])
17
```



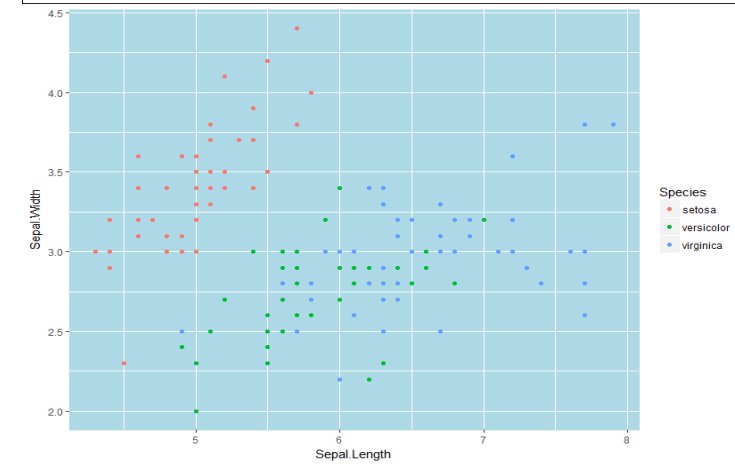
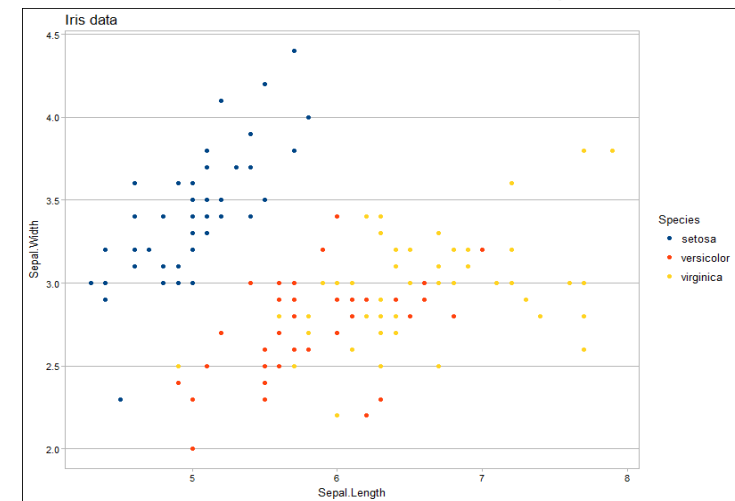
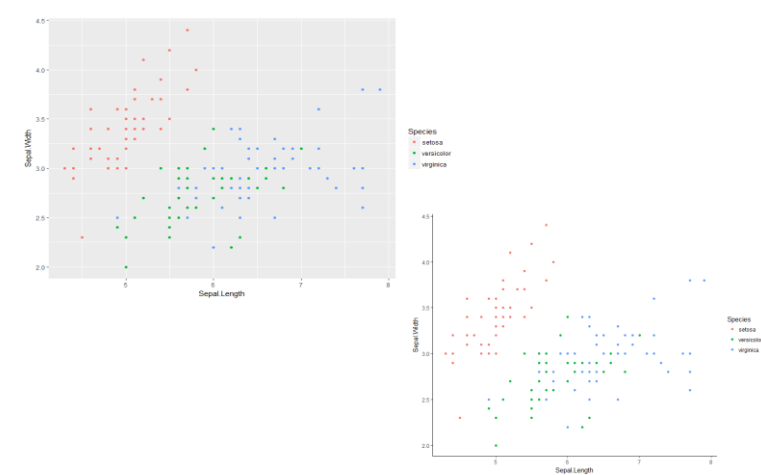
气泡图 (Bubbles)

```
1 require(ggplot2)
2 df<- read.csv("Bubbles.csv")
3 > df
4   Term Input_number P.Value id
5   1 Ribosome 45 0.0000000000000000797516 1
6   2 Biosynthesis of secondary metabolites 354 0.000000000000000149654 2
7   3 Metabolic pathways 282 0.000000000000394082733 3
8   4 Plant hormone signal transduction 193 0.00000001163614774940 4
9   5 Galactose metabolism 29 0.00001489712860000000 5
10  6 Flavonoid biosynthesis 18 0.00048351702379700000 6
11  7 Carbon fixation in photosynthetic organisms 30 0.00000197712300000000 7
12  8 Fatty acid elongation 20 0.00000852890000000000 8
13  9 Fatty acid metabolism 424 0.0002962069999999999999 9
14 10 Fatty acid degradation 58 0.0003399950000000000000 10
15 11 Steroid biosynthesis 16 0.0005169999999999999999 11
16 12 Valine, leucine and isoleucine degradation 19 0.0005875519999999999997 12
17 13 Phagosome 67 0.0005879299999999999996 13
18 14 Starch and sucrose metabolism 48 0.0044200000000000000032 14
19 15 Biosynthesis of unsaturated fatty acids 23 0.0025335000000000000016 15
20
21 ggplot(df, aes(x = id,y=Term,label = Term)) +
22   geom_point(aes(size = Input_number, colour = P.Value)) +
23   #geom_text(hjust = 1, size = 2) +
24   scale_size(range = c(1,15)) +
25   scale_x_continuous(breaks = seq(1, 15, 2)) +
26   scale_colour_gradientn(colours=rainbow(4)) +
27   theme_bw()
28
```



美化 (themes and background)

```
1 # ggplot2自带主题
2 p <- ggplot(iris, aes(Sepal.Length, Sepal.Width, colour = Species))+
3   geom_point()
4 p
5 p + theme_classic()
6 -----
7 # 主题包
8 install.packages("ggthemes") # Install
9 library(ggthemes) # Load
10 p + theme_calc()+ scale_colour_calc()+
11   ggtitle("Iris data")
12 -----
13 # 定制主题
14 p + theme(
15   panel.background = element_rect(fill = "lightblue",
16                                   colour = "lightblue",
17                                   size = 0.5, linetype = "solid"),
18   panel.grid.major = element_line(size = 0.5, linetype = 'solid',
19                                   colour = "white"),
20   panel.grid.minor = element_line(size = 0.25, linetype = 'solid',
21                                   colour = "white")
22 )
23
```

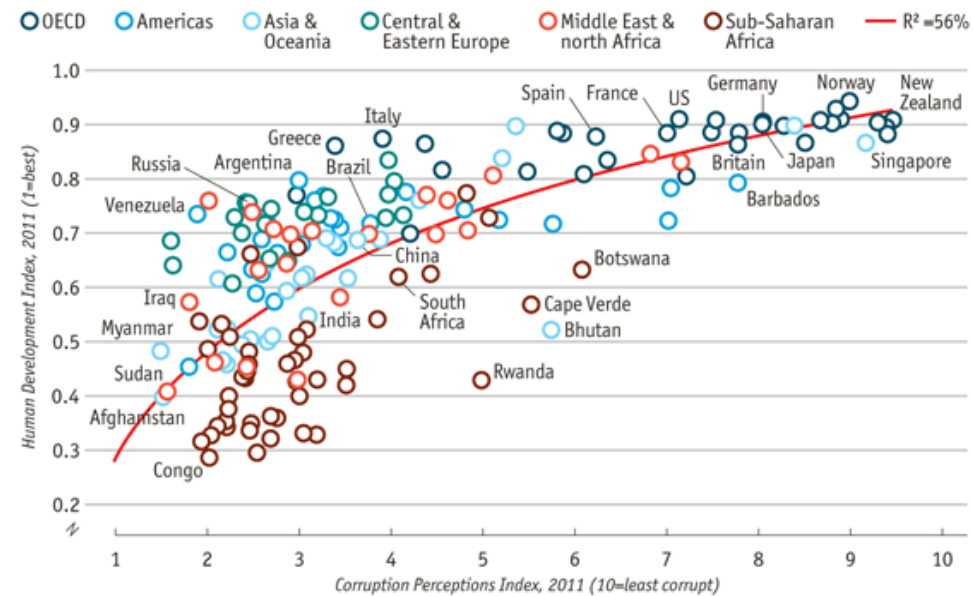


RNA-Seq (DESeq2)

```
1 library(DESeq2);library(limma);library(pasilla);data(pasillaGenes);exprSet=counts(pasillaGenes)
2 # 表达矩阵
3 head(exprSet)
4   treated1fb treated2fb treated3fb untreated1fb untreated2fb untreated3fb untreated4fb
5 FBgn0000003          0          0          1          0          0          0          0
6 FBgn0000008          78          46          43          47          89          53          27
7 FBgn0000014          2          0          0          0          0          1          0
8 FBgn0000015          1          0          1          0          1          1          2
9 FBgn0000017        3187        1672        1859        2445        4615        2063        1711
10 FBgn0000018         369         150         176         288         383         135         174
11 # 构建dds对象
12 colData <- data.frame(row.names=colnames(exprSet), group_list=group_list)
13 > colData
14   group_list
15 treated1fb treated
16 treated2fb treated
17 treated3fb treated
18 untreated1fb untreated
19 untreated2fb untreated
20 untreated3fb untreated
21 untreated4fb untreated
22 dds <- DESeqDataSetFromMatrix(countData = exprSet,colData = colData,design = ~ group_list)
23 # normalization
24 dds2 <- DESeq(dds)
25 # 提取差异分析结果
26 resultsNames(dds2)
27 res <- results(dds2, contrast=c("group_list","treated","untreated"))
28 resOrdered <- res[order(res$padj),]
29 resOrdered=as.data.frame(resOrdered)
30 head(resOrdered)
31   baseMean log2FoldChange lfcSE stat pvalue padj
32 FBgn0039155 453.2753 -3.714214 0.1600580 -23.20543 4.013291e-119 3.089431e-115
33 FBgn0029167 2165.0445 -2.082793 0.1035963 -20.10491 6.684454e-90 2.572846e-86
34 FBgn0035085 366.8279 -2.227243 0.1369744 -16.26028 1.888618e-59 4.846194e-56
35 FBgn0029896 257.9027 -2.206780 0.1586969 -13.90563 5.854593e-44 1.126716e-40
36 FBgn0034736 118.4074 -2.565002 0.1847628 -13.88268 8.067448e-44 1.242064e-40
37 FBgn0040091 610.6035 -1.430433 0.1201539 -11.90501 1.114552e-32 1.429970e-29
38
```


复杂图形修改

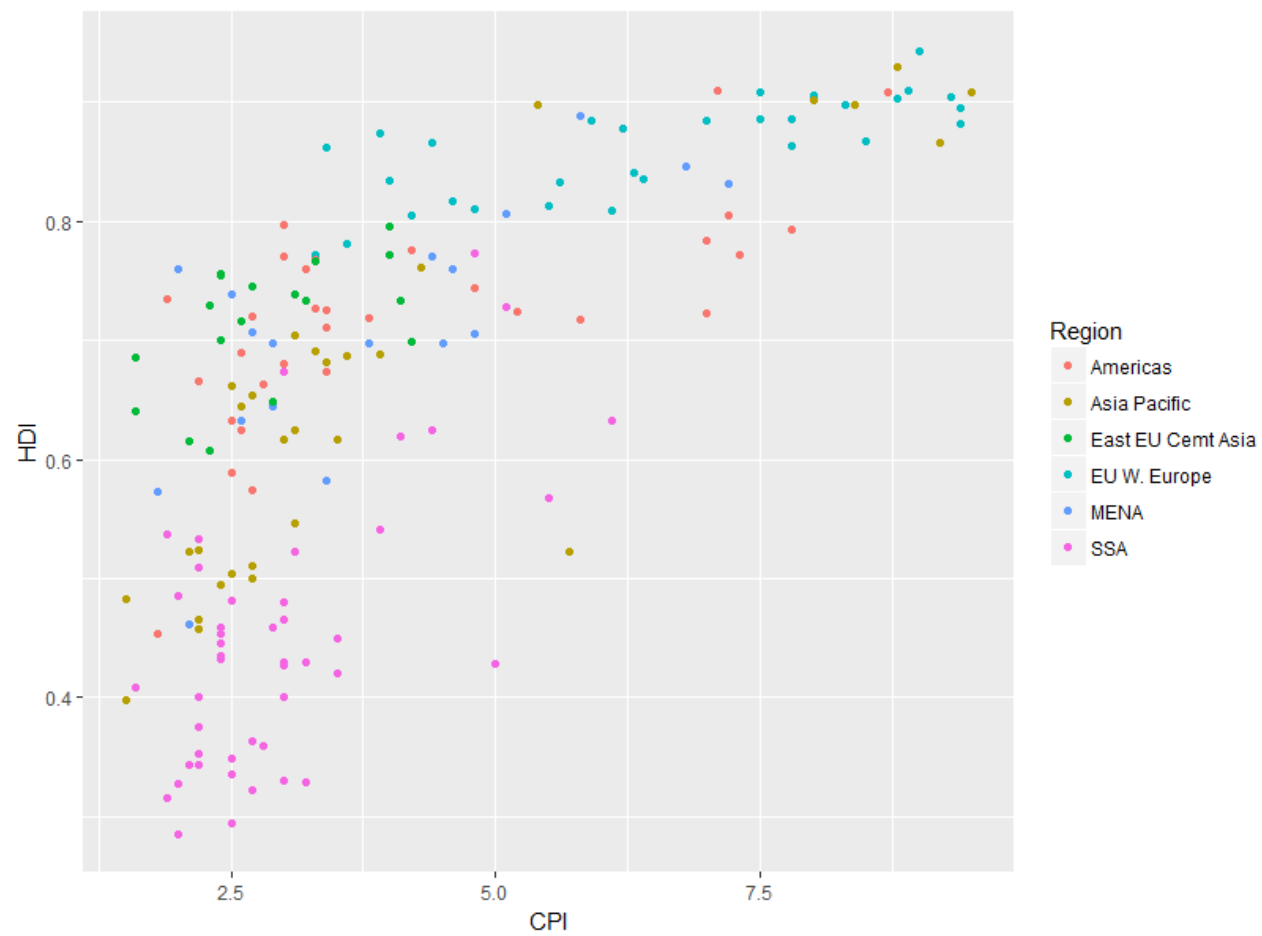
Corruption and human development



Sources: Transparency International; UN Human Development Report

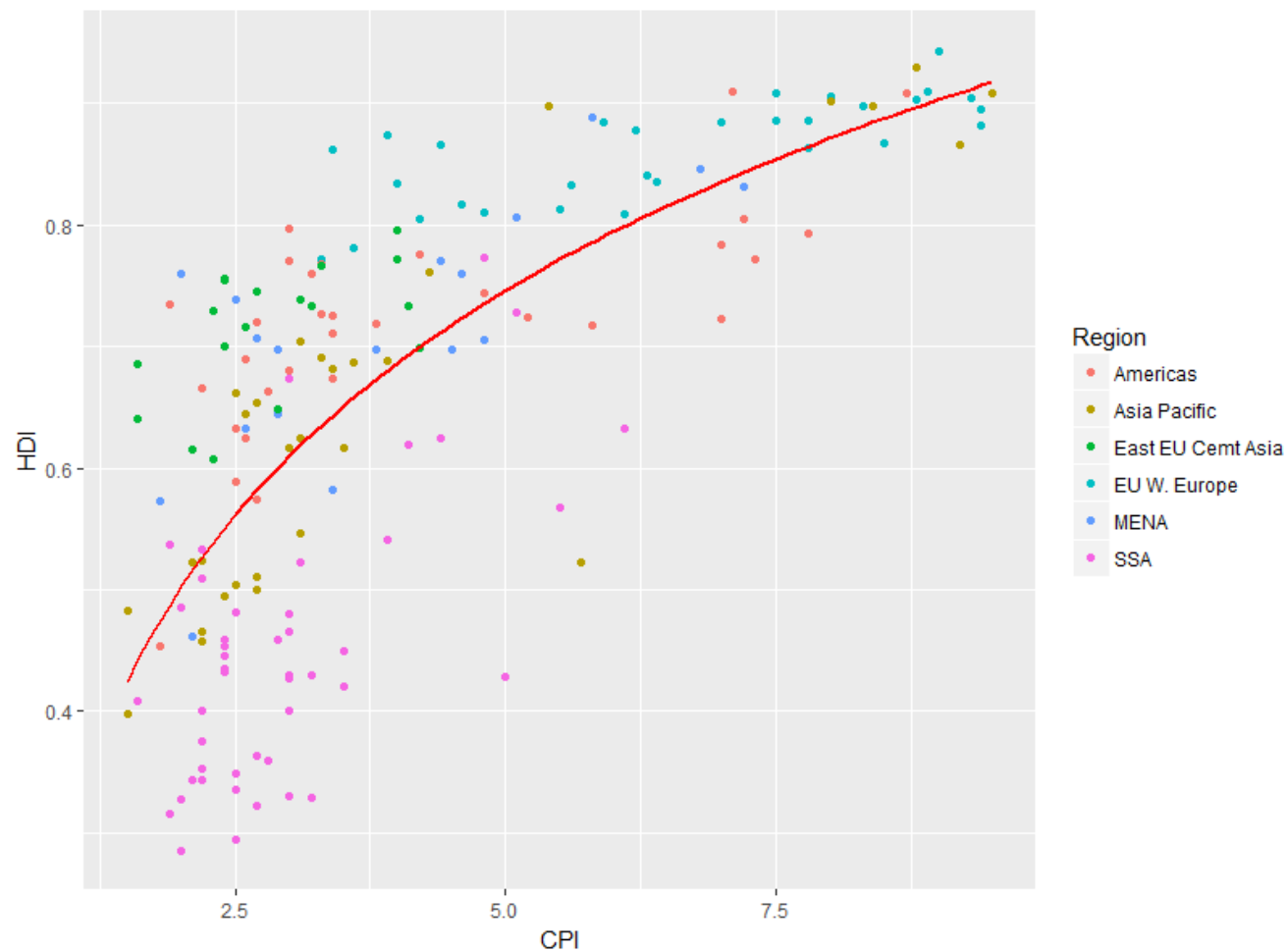
1. 基本图形

```
1 library(ggplot2)
2 dat <- read.csv("EconomistData.csv")
3 # Basic plot
4 pc1 <- ggplot(dat, aes(x = CPI, y = HDI, color = Region)) +
5   geom_point()
6 pc1
7
```



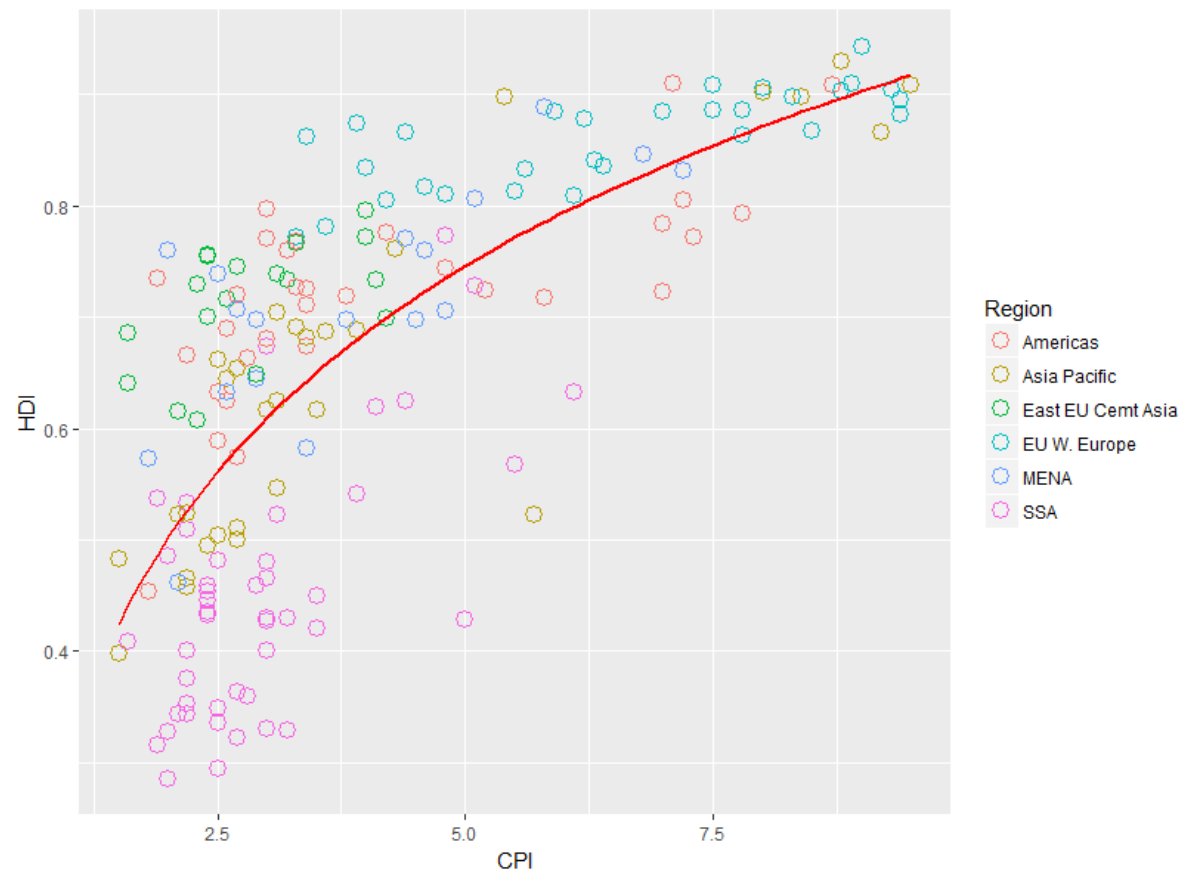
2. 添加趋势线

```
1 # Trend line
2 pc2 <- pc1 +
3   geom_smooth(aes(group = 1),
4               method = "lm",
5               formula = y ~ log(x),
6               se = FALSE,
7               color = "red")
8 pc2
9
```



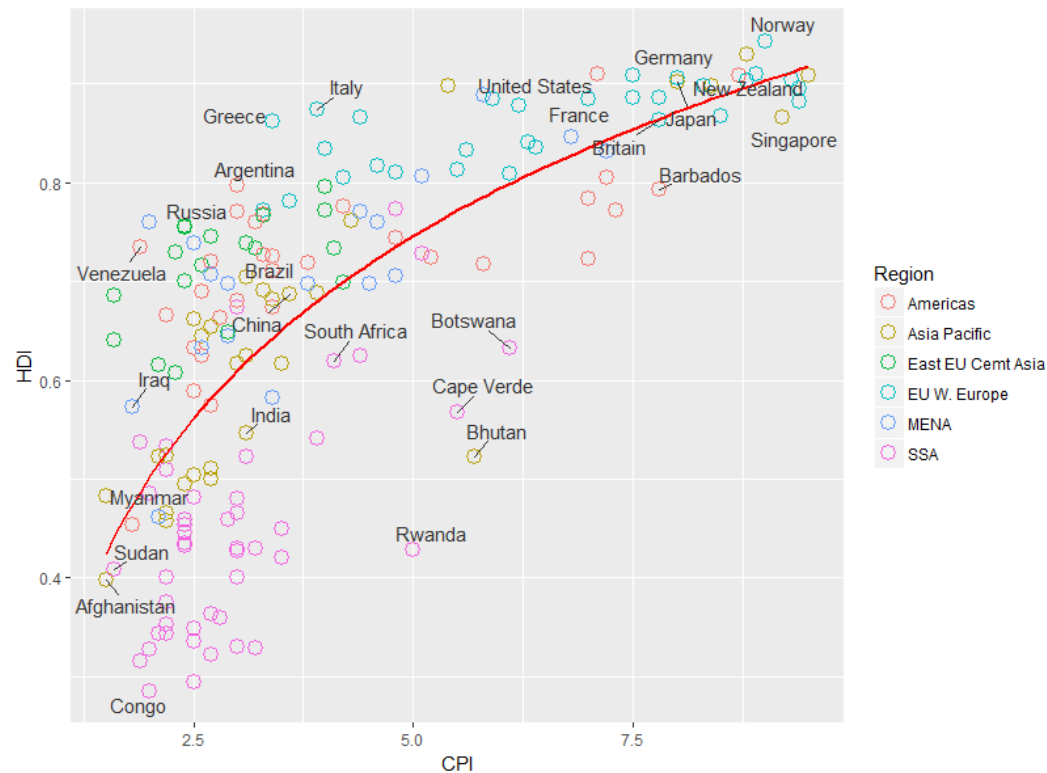
3. 圆卷

```
1  
2  
3 # Open points  
4 pc3 <- ggplot(dat, aes(x = CPI, y = HDI, color = Region))+  
5   geom_point(shape = 1, size = 4) +  
6   geom_smooth(aes(group = 1),  
7               method = "lm",  
8               formula = y ~ log(x),  
9               se = FALSE,  
10              color = "red")  
11 pc3  
12
```



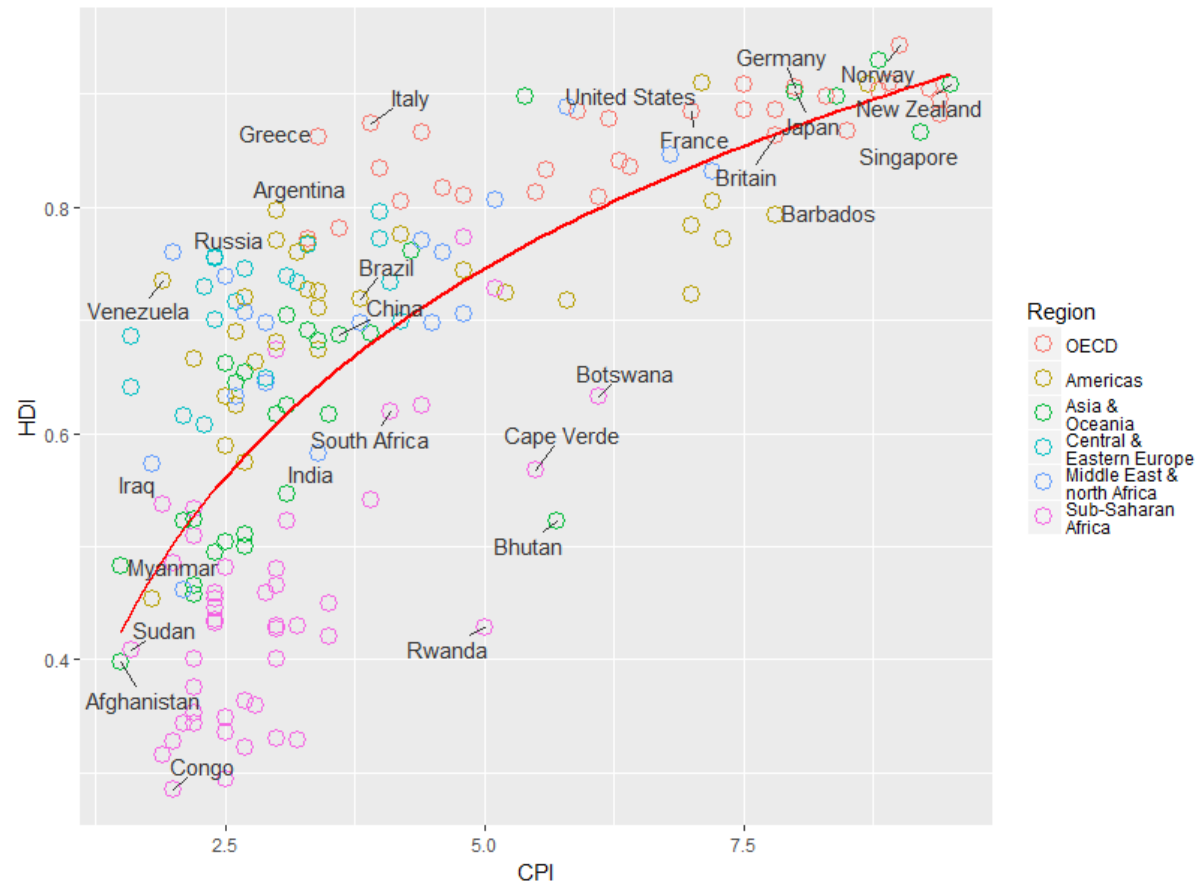
4. 标注想要的点

```
1 pointsToLabel <- c("Russia", "Venezuela", "Iraq", "Myanmar", "Sudan",  
2 "Afghanistan", "Congo", "Greece", "Argentina", "Brazil",  
3 "India", "Italy", "China", "South Africa", "Spain",  
4 "Botswana", "Cape Verde", "Bhutan", "Rwanda",  
5 "France", "United States", "Germany", "Britain", "Barbados",  
6 "Norway", "Japan",  
7 "New Zealand", "Singapore")  
8 library("ggrepel")  
9 pc4 <- pc3 + geom_text_repel(aes(label = Country),  
10 color = "gray20",  
11 data = subset(dat, Country %in% pointsToLabel),  
12 force = 10)  
13 pc4  
14
```



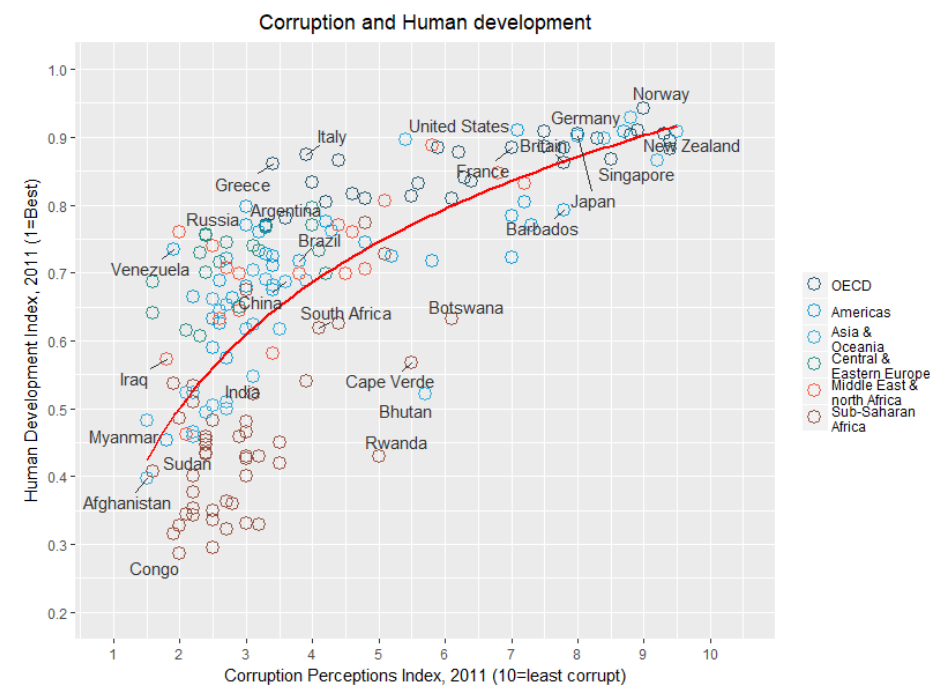
5. 修改图例值和顺序

```
1 dat$Region <- factor(dat$Region,  
2                       levels = c("EU W. Europe",  
3                                   "Americas",  
4                                   "Asia Pacific",  
5                                   "East EU Cemt Asia",  
6                                   "MENA",  
7                                   "SSA"),  
8                       labels = c("OECD",  
9                                   "Americas",  
10                                  "Asia &\nOceania",  
11                                  "Central &\nEastern Europe",  
12                                  "Middle East &\nnorth Africa",  
13                                  "Sub-Saharan\nAfrica"))  
14 pc4$data <- dat  
15 pc4  
16
```



6. 利用scale来修改x, y轴, 颜色和标出title

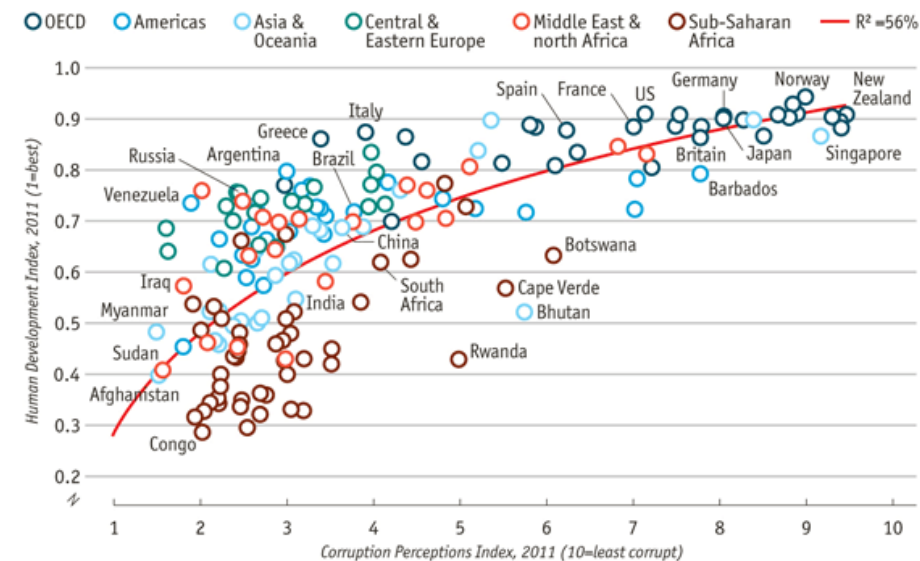
```
1 library(grid)
2 pc5 <- pc4 +
3   scale_x_continuous(name = "Corruption Perceptions Index, 2011 (10=least corrupt)",
4                     limits = c(.9, 10.5),
5                     breaks = 1:10) +
6   scale_y_continuous(name = "Human Development Index, 2011 (1=Best)",
7                     limits = c(0.2, 1.0),
8                     breaks = seq(0.2, 1.0, by = 0.1)) +
9   scale_color_manual(name = "",
10                    values = c("#24576D",
11                              "#099DD7",
12                              "#28AADC",
13                              "#248E84",
14                              "#F2583F",
15                              "#96503F")) +
16   ggtitle("Corruption and Human development") +
17   theme(plot.title = element_text(hjust = 0.5))
18 pc5
19
```



7. 微调主题

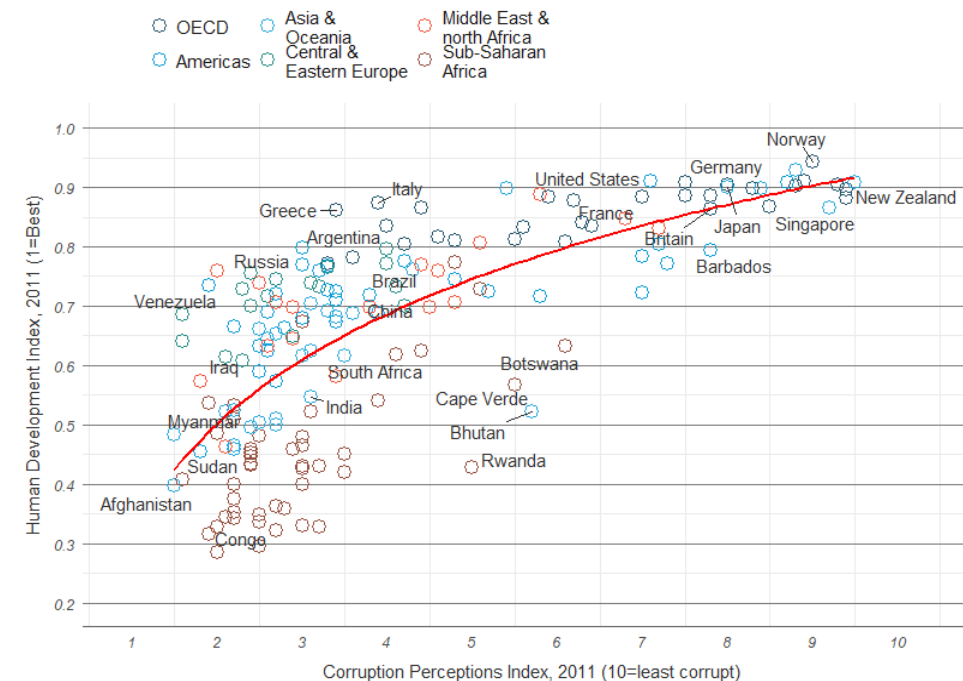
```
1 library(grid) # for the 'unit' function
2 pc6 <- pc5 +
3   theme_minimal() + # start with a minimal theme and add what we need
4   theme(text = element_text(color = "gray20"),
5         legend.position = c("top"), # position the legend in the upper left
6         legend.direction = "horizontal",
7         legend.justification = 0.1, # anchor point for legend.position.
8         legend.text = element_text(size = 11, color = "gray10"),
9         axis.text = element_text(face = "italic"),
10        axis.title.x = element_text(vjust = -1), # move title away from axis
11        axis.title.y = element_text(vjust = 2), # move away for axis
12        axis.ticks.y = element_blank(), # element_blank() is how we remove elements
13        axis.line = element_line(color = "gray40", size = 0.5),
14        axis.line.y = element_blank(),
15        panel.grid.major = element_line(color = "gray50", size = 0.5),
16        panel.grid.major.x = element_blank()
17      )
18 pc6
19
```

Corruption and human development



Sources: Transparency International; UN Human Development Report

Corruption and Human development



THANKS FOR YOUR WATCHING



<http://tiramisutes.github.io/>

@Huazhong Agricultural University

