

[Home](#)[About](#)[Categories](#)[Tags](#)

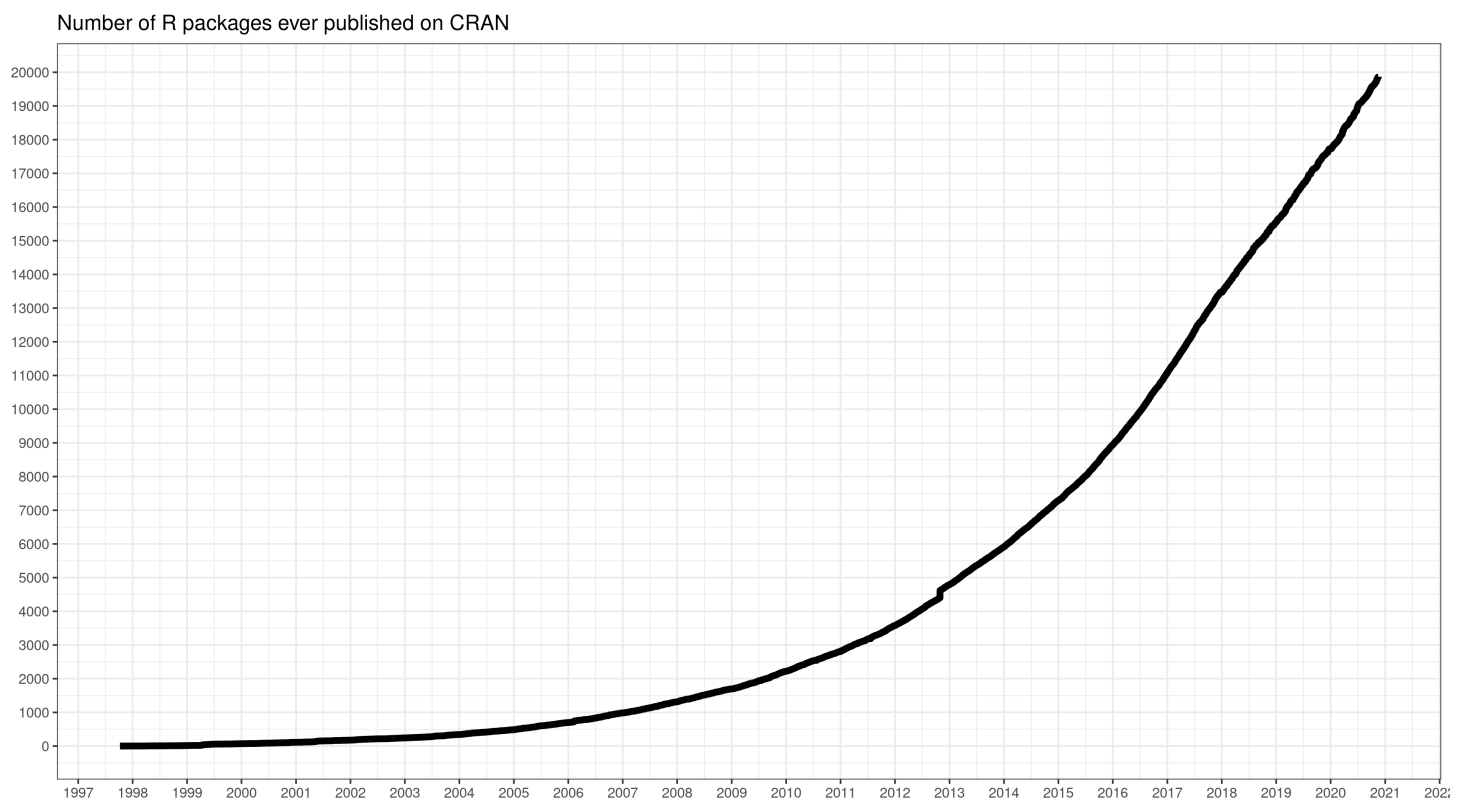
Learning R

Pylone

2021/08/05

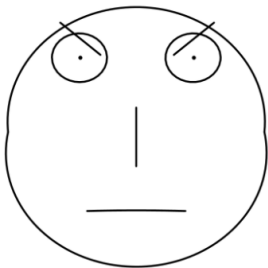
1. Why

1.1 开源：相当于手机系统，可以使用任何在此平台上开发的软件，叫做包（**package**）

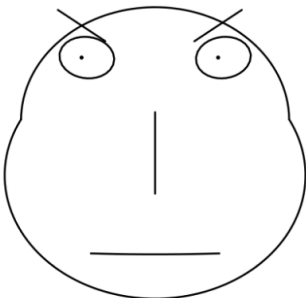


packages number

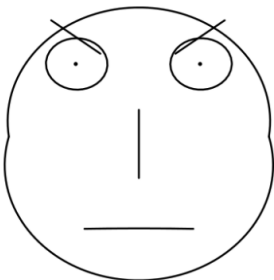
1.2 无所不能的分析及可视化工具，一些可视化例子，浏览以便于激发灵感



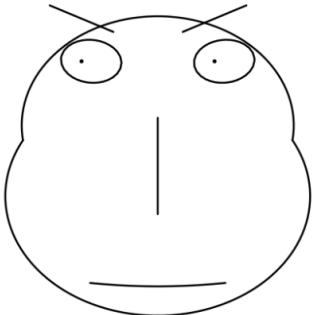
AARONSON,L.H.



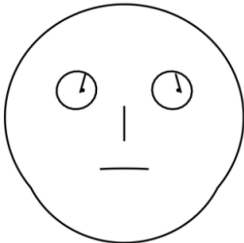
ALEXANDER,J.M.



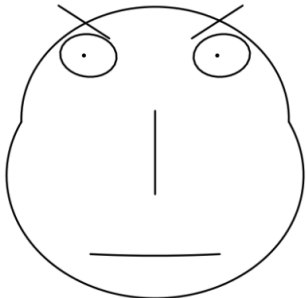
ARMENTANO,A.J.



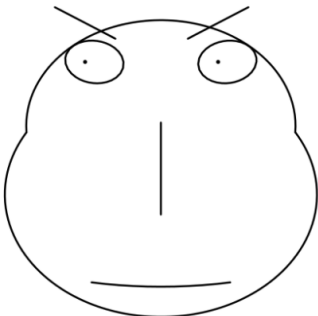
BERDON,R.I.



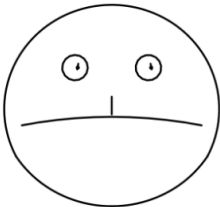
BRACKEN,J.J.



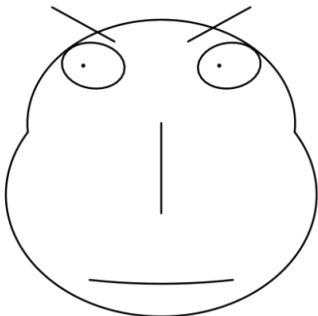
BURNS,E.B.



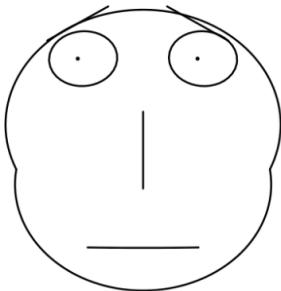
CALLAHAN,R.J.



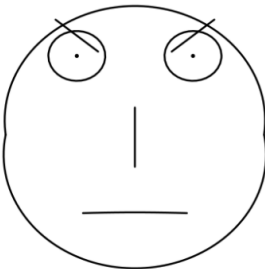
COHEN,S.S.



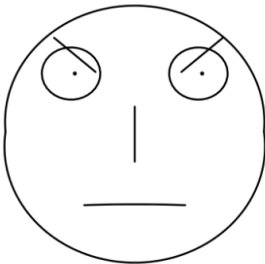
DALY,J.J.



DANNEHY,J.F.



DEAN,H.H.



DEVITA,H.J.

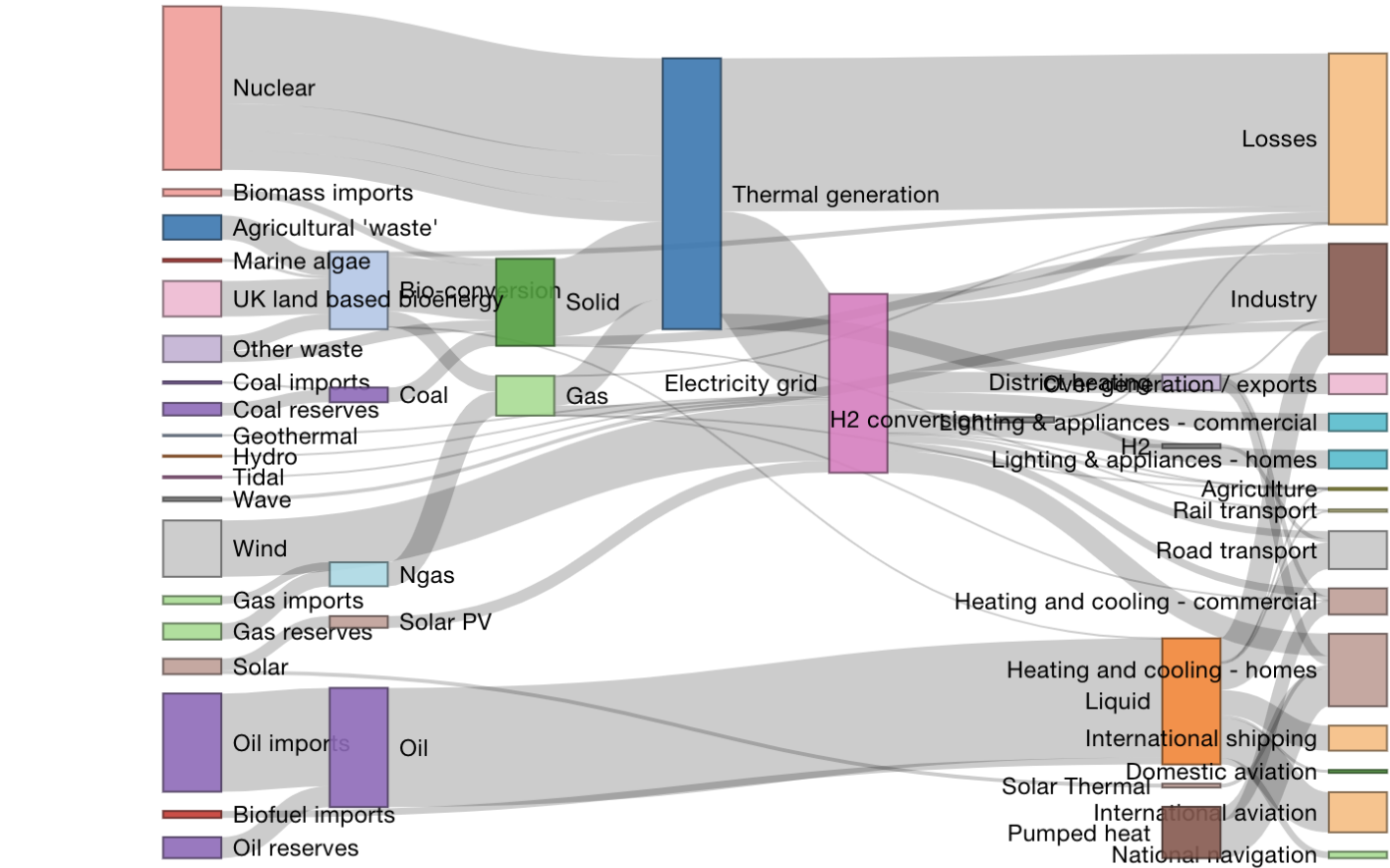
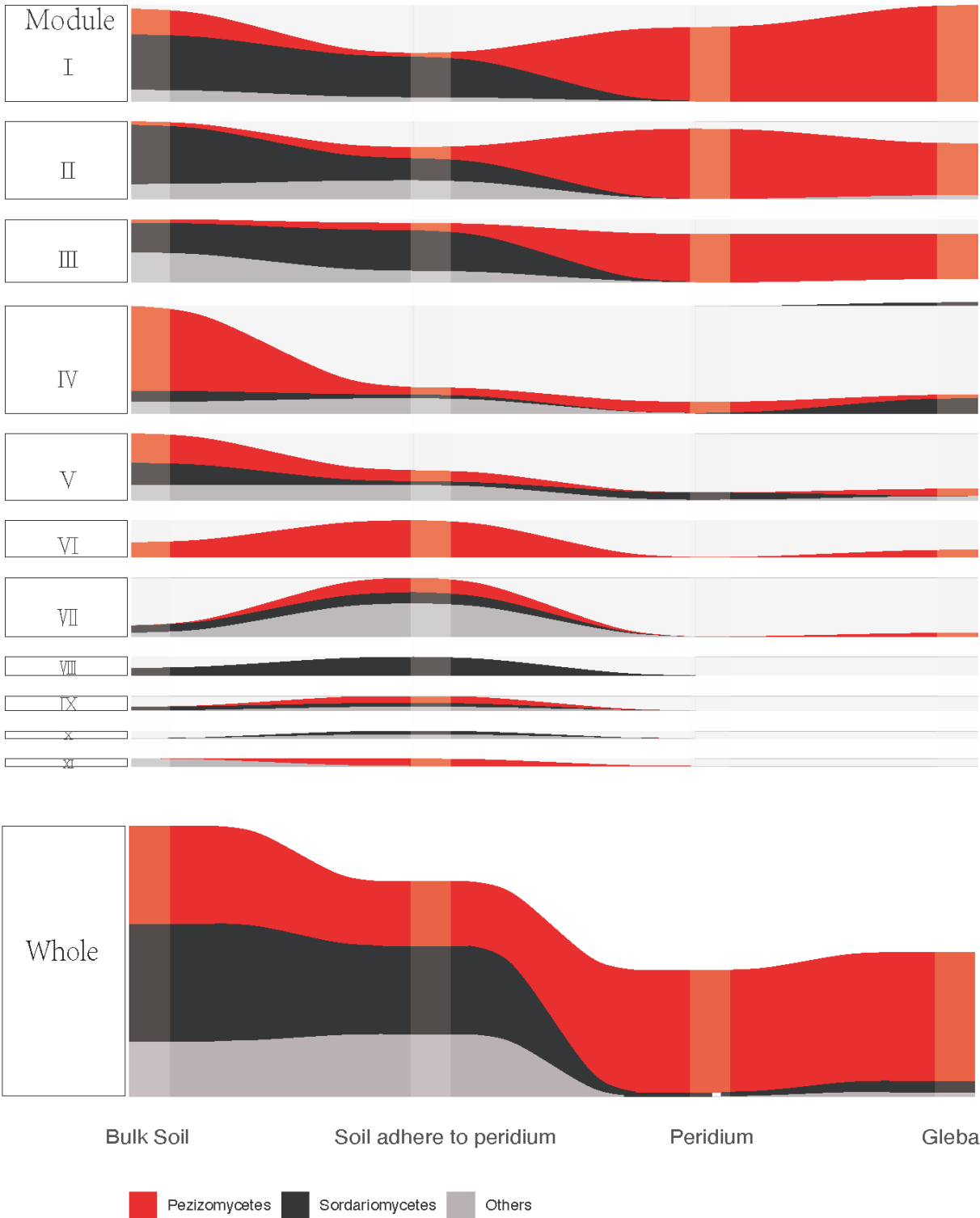
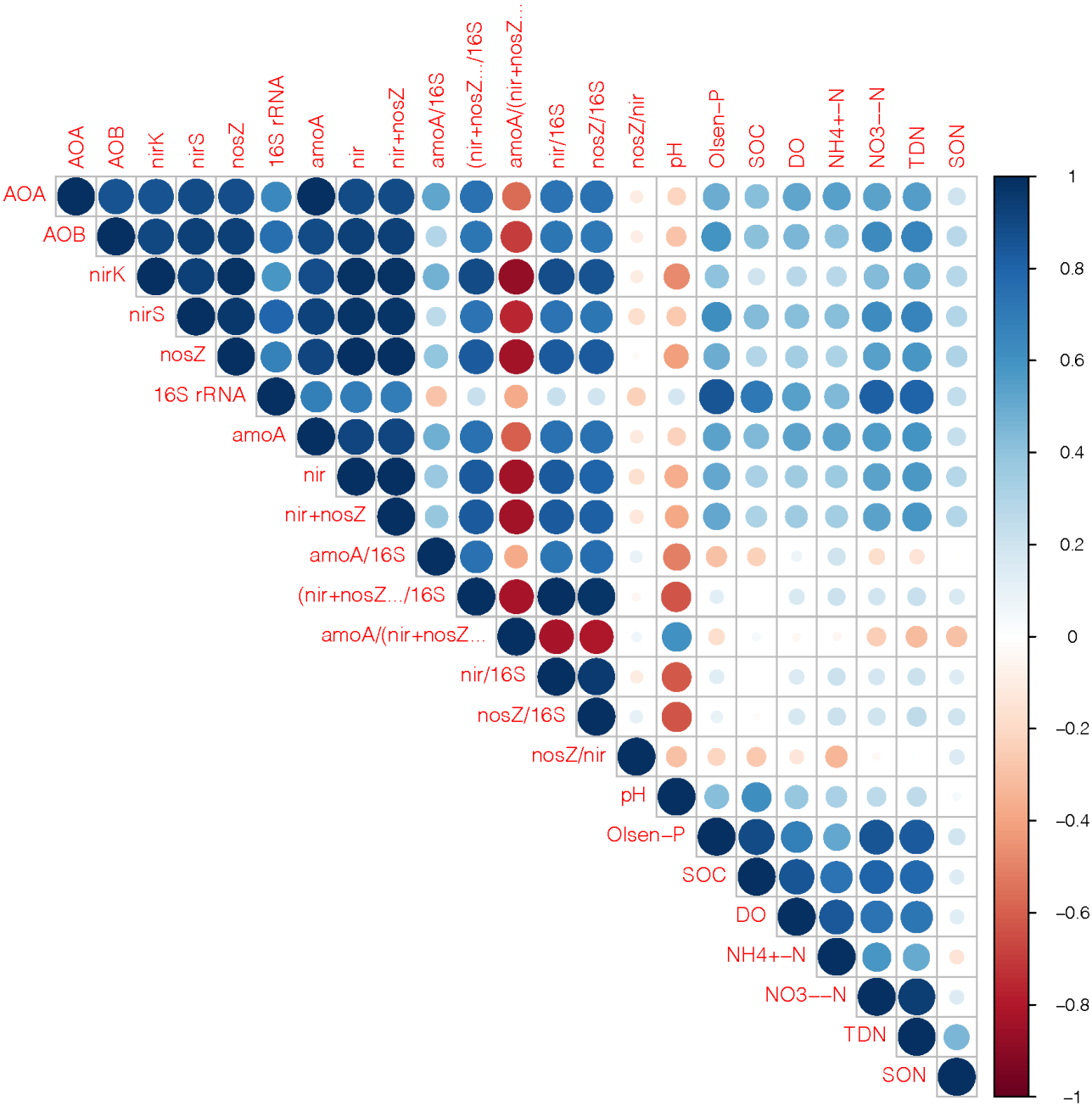


Figure 9.8: Sankey diagram

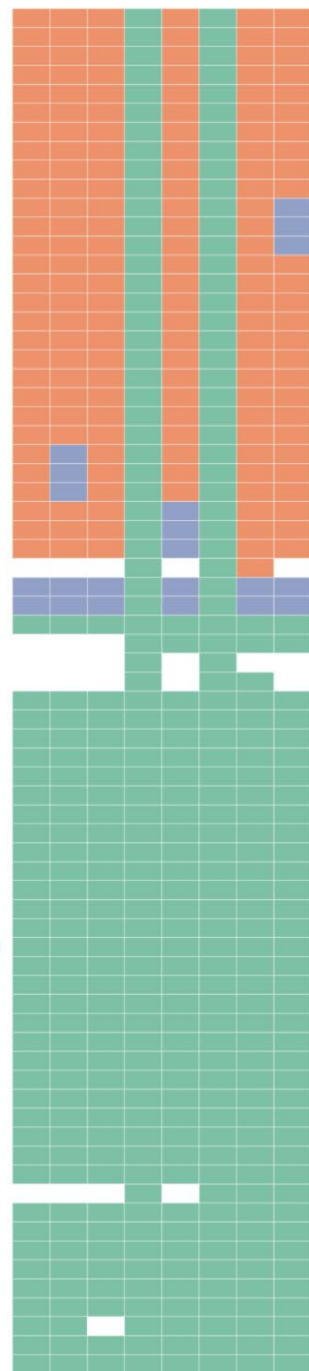
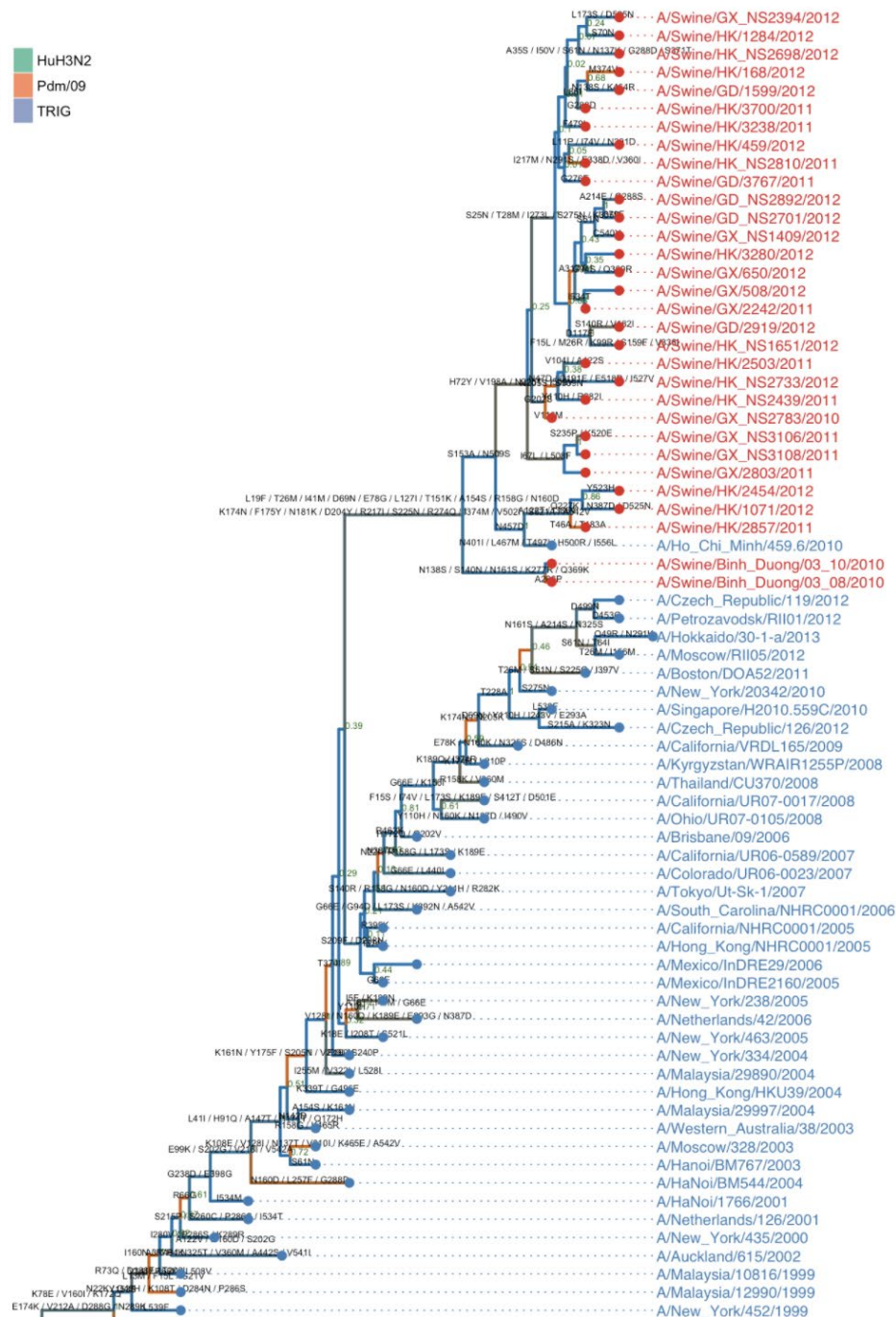
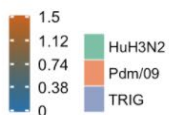
桑基图



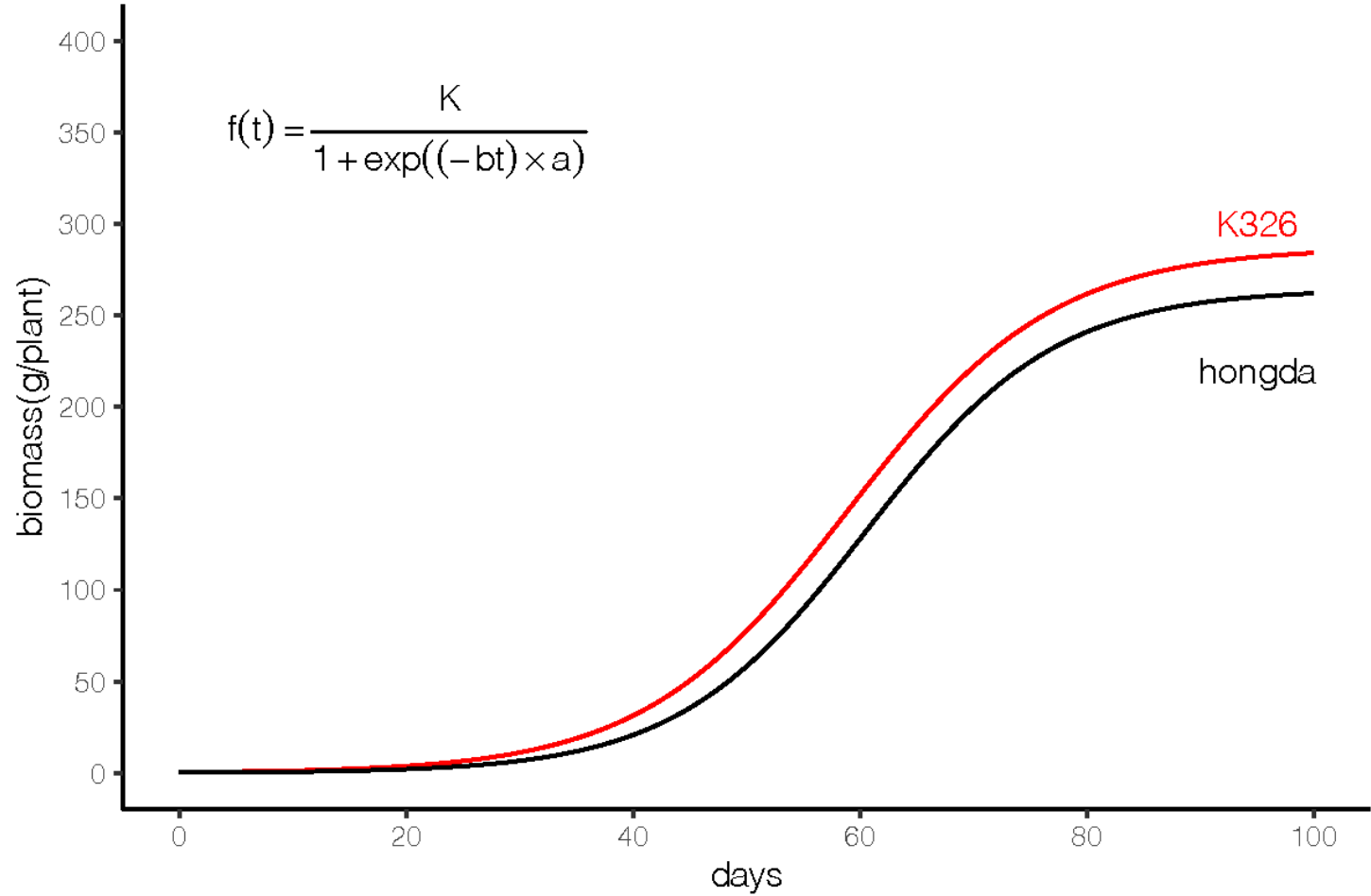
冲积图



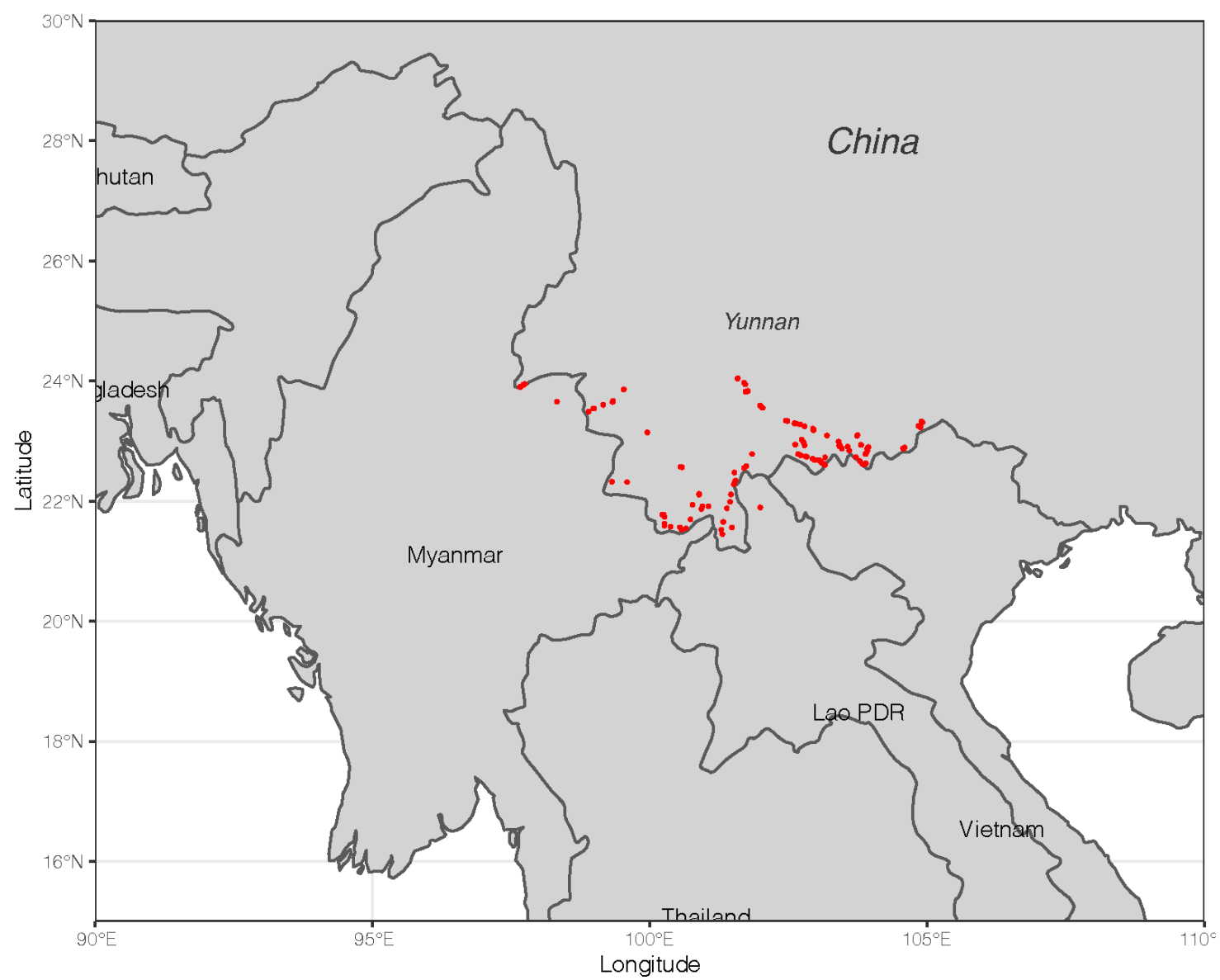
相关性



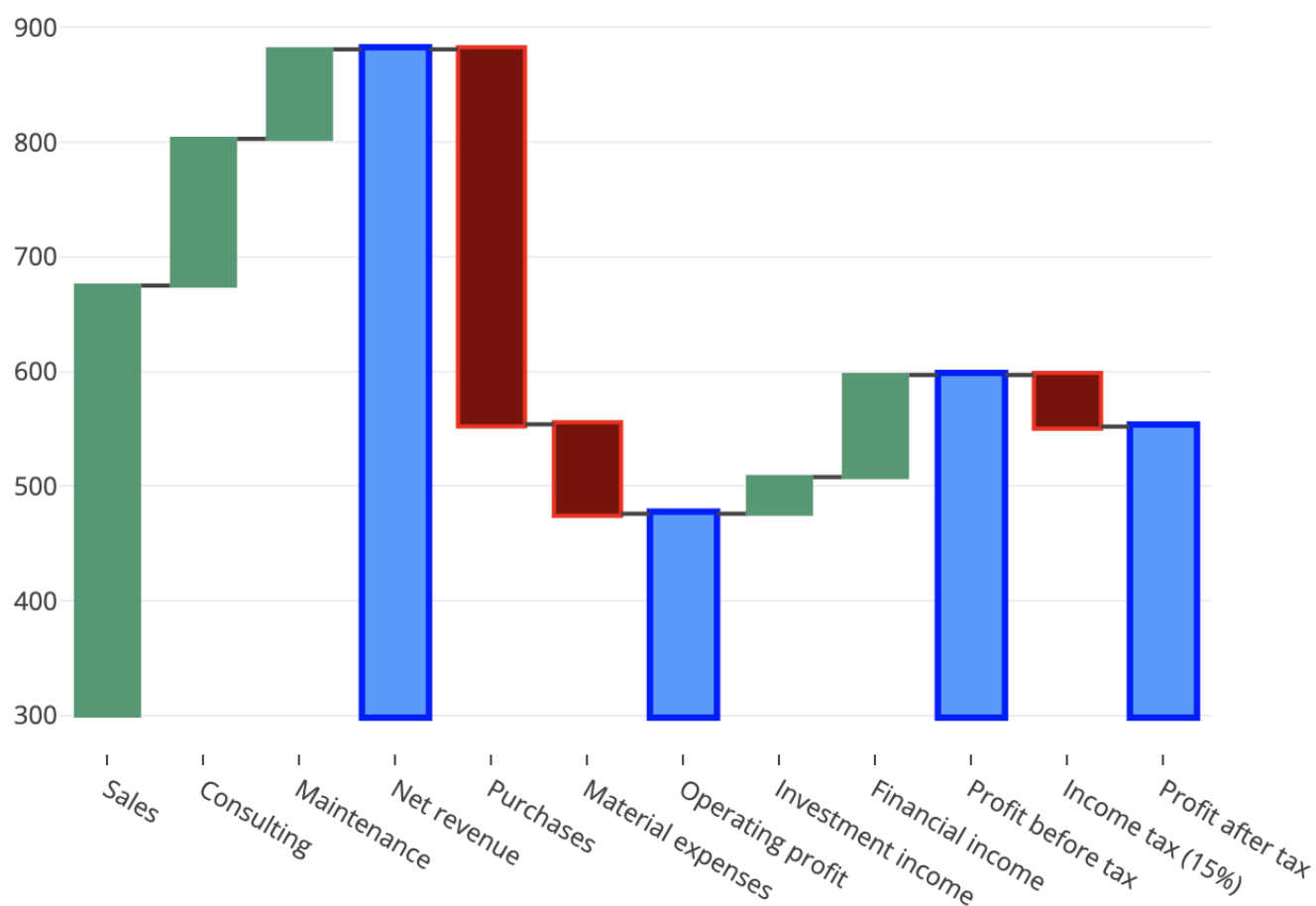
聚类图



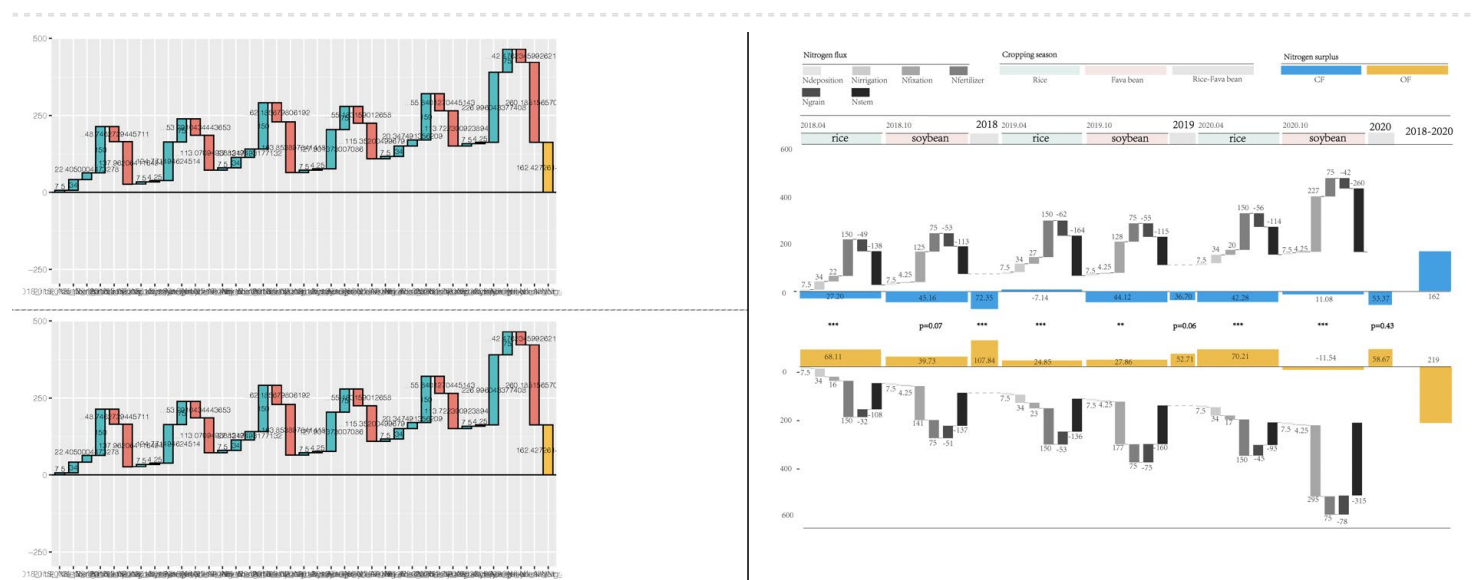
曲线图



map



waterfall



使用R出图，AI加工图

1.3 更大魅力在于对数据的自由探索

1.4丰富的功能：word、pdf、ppt、html、blog等

2.How

Q1:怎么学的更快？

带着需求来学习，学以致用

Q2:小白怎么入手？

安装、基础操作：教程，《R语言实战》相当于查询手册

Q3:遇到问题怎么办？

(1) 先从教程（查询手册）中找答案，再去百度、谷歌（推荐），再交流、总结、分享

(2) 真想学好还是建议翻墙，巨量优质的教程分享

ggplot2数据可视化::速查表



基础

ggplot2 基于图形语法，使用相同的组件（数据集、坐标系统和表示数据点的几何对象）来构建图片。



为了获取显示值，数据中的变量映射到图形的视觉属性，如大小、颜色以及x和y位置。



完成以下模板来构建图形

```
ggplot(data = <DATA>) +
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),
    stat = <STAT>, position = <POSITION>) +
    <COORDINATE_FUNCTION> +
    <FACET_FUNCTION> +
    <SCALE_FUNCTION> +
    <THEME_FUNCTION>
```

ggplot(data = mpg, aes(x = cty, y = hwy)) 通过添加图层来完成图形，每层添加一个geom函数。

qplot(x = cty, y = hwy, data = mpg, geom = "point") 用给定的数据、几何对象和映射创建完整的图片。绘图函数提供许多有用的默认设置。

last_plot() 返回上一个图片

ggsave("plot.png", width = 5, height = 5) 将最后一个图片保存至工作目录中名为"plot.png"的5'x'5'文件。文件类型与文件扩展名相匹配。

几何对象

使用geom函数表示数据点，使用geom的属性表示变量。每个函数绘制一个图层。

基本图像

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))

a + geom_blank()
(Useful for expanding limits)

b + geom_curve(aes(yend = lat + 1,
  xend = long - 1, curvature = z)) ~ x, yend, y, alpha,
  angle, color, curvature, linetype, size

a + geom_path(lineend = "butt", linejoin = "round",
  linemitre = 1)
x, y, alpha, color, group, linetype, size

a + geom_polygon(aes(group = group))
x, y, alpha, color, fill, group, linetype, size

b + geom_rect(aes(xmin = long, ymin = lat, xmax = long
  + 1, ymax = lat + 1)) ~ x, ymax, xmin, ymax, ymin, alpha,
  color, fill, linetype, size

a + geom_ribbon(aes(ymin = unemploy - 900,
  ymax = unemploy + 900)) ~ x, ymax, ymin, alpha,
  color, fill, group, linetype, size
```

分段线

```
常用参数: x, y, alpha, color, linetype, size
b + geom_abline(aes(intercept = 0, slope = 1))
b + geom_hline(aes(yintercept = lat))
b + geom_vline(aes(xintercept = long))

b + geom_segment(aes(yend = lat + 1, xend = long + 1))
b + geom_spoke(aes(angle = 1:115, radius = 1))
```

单一变量 连续

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

c + geom_area(stat = "bin")
x, y, alpha, color, fill, linetype, size

c + geom_density(kernel = "gaussian")
x, y, alpha, color, fill, group, linetype, size, weight

c + geom_dotplot()
x, y, alpha, color, fill

c + geom_freqpoly() x, y, alpha, color, group, linetype, size

c + geom_histogram(binwidth = 5) x, y, alpha, color, fill, linetype, size, weight

c2 + geom_qq(aes(sample = hwy)) x, y, alpha, color, fill, linetype, size, weight
```

离散

```
d <- ggplot(mpg, aes(fll))

d + geom_bar()
x, alpha, color, fill, linetype, size, weight
```

双变量

```
连续x、连续y
e <- ggplot(mpg, aes(cty, hwy))

e + geom_label(aes(label = cty), nudge_x = 1,
  nudge_y = 1, check_overlap = TRUE) x, y, label,
  alpha, angle, color, family, fontface, hjust,
  lineheight, size, vjust

e + geom_jitter(height = 2, width = 2)
x, y, alpha, color, fill, shape, size

e + geom_point(), x, y, alpha, color, fill, shape,
  size, stroke

e + geom_quantile(), x, y, alpha, color, group,
  linetype, size, weight

e + geom_rug(sides = "bl"), x, y, alpha, color,
  linetype, size

e + geom_smooth(method = lm), x, y, alpha,
  color, fill, group, linetype, size, weight

e + geom_text(aes(label = cty), nudge_x = 1,
  nudge_y = 1, check_overlap = TRUE) x, y, label,
  alpha, angle, color, family, fontface, hjust,
  lineheight, size, vjust
```

离散x、连续y

```
f <- ggplot(mpg, aes(class, hwy))

f + geom_col(), x, y, alpha, color, fill, group,
  linetype, size

f + geom_boxplot(), x, y, lower, middle, upper,
  ymax, ymin, alpha, color, fill, group, linetype,
  shape, size, weight

f + geom_dotplot(binaxis = "y", stackdir =
  "center"), x, y, alpha, color, fill, group

f + geom_violin(scale = "area"), x, y, alpha,
  color, fill, group, linetype, size, weight
```

离散x、离散y

```
g <- ggplot(diamonds, aes(cut, color))

g + geom_count(), x, y, alpha, color, fill, shape,
  size, stroke
```

三变量

```
sealsSz <- with(seals, sqrt(delta_long^2 + delta_lat^2))
l <- ggplot(seals, aes(long, lat))

l + geom_contour(aes(z = z))
x, y, z, alpha, colour, group, linetype,
  size, weight

l + geom_raster(aes(fill = z), hjust=0.5, vjust=0.5,
  interpolate=FALSE)
x, y, alpha, fill

l + geom_tile(aes(fill = z)), x, y, alpha, color, fill,
  linetype, size, width
```

连续二元分布

```
h <- ggplot(diamonds, aes(carat, price))

h + geom_bin2d(binwidth = c(0.25, 500))
x, y, alpha, color, fill, linetype, size, weight

h + geom_density2d()
x, y, alpha, colour, group, linetype, size

h + geom_hex()
x, y, alpha, colour, fill, size
```

连续函数

```
i <- ggplot(economics, aes(date, unemploy))

i + geom_area()
x, y, alpha, color, fill, linetype, size

i + geom_line()
x, y, alpha, colour, group, linetype, size

i + geom_step(direction = "hv")
x, y, alpha, color, group, linetype, size
```

误差的呈现方式

```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))

j + geom_crossbar(fatten = 2)
x, y, ymax, ymin, alpha, color, fill, group, linetype, size

j + geom_errorbar(), x, ymax, ymin, alpha, color,
  group, linetype, size, width (also geom_errorbarh())

j + geom_linerange()
x, ymin, ymax, alpha, color, group, linetype, size

j + geom_pointrange()
x, y, ymin, ymax, alpha, color, fill, group, linetype,
  shape, size
```

地图

```
data <- data.frame(murder = USArrests$Murder,
  state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))

k + geom_map(aes(map_id = state), map = map)
+ expand_limits(x = map$long, y = map$lat),
  map_id, alpha, color, fill, linetype, size
```



RStudio® is a trademark of RStudio, Inc. • [CC BY SA](#) RStudio • [info@rstudio.com](#) • 844-448-1212 • [rstudio.com](#) • Learn more at <http://ggplot2.tidyverse.org> • ggplot2 2.1.0 • Updated: 2016-11

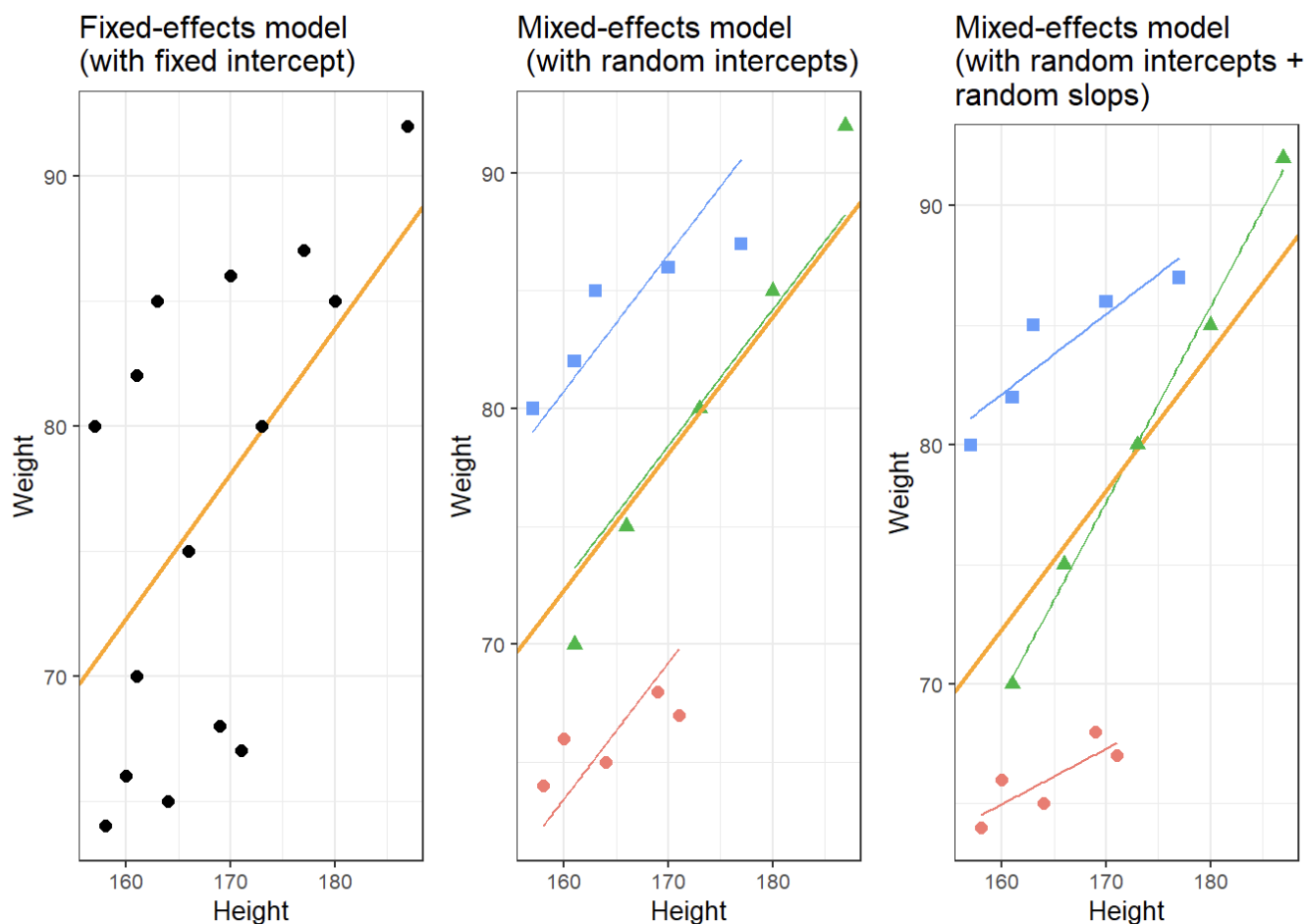
速查表

3.我的数据可以用哪些方法分析？多元统计

逃不出多元统计分析，去看学习强国的多元统计分析课程。

3.1 变量

- 自变量：多元相关cor、主成分分析OCA、因子分析CA
- 自变量与因变量：多元回归、线性模型（广义glm、一般（方差分析）、普通、线性混合效应模型lme）、冗余分析RDA



线性混合效应模型

尤其适用于以下情况：

- 数据有缺失值；
- 不平衡设计，如一个处理3个样本，另一个5个样本；
- 结构化（nested）数据；
- 结果变量非连续，比如评分（5，10，15）；
- 等等。。

3.2样本

- 聚类分析

3.3样本与变量

- 样本与自变量：对应分析（CA）
- 样本、自变量与因变量：约束性对应分析（CCA）

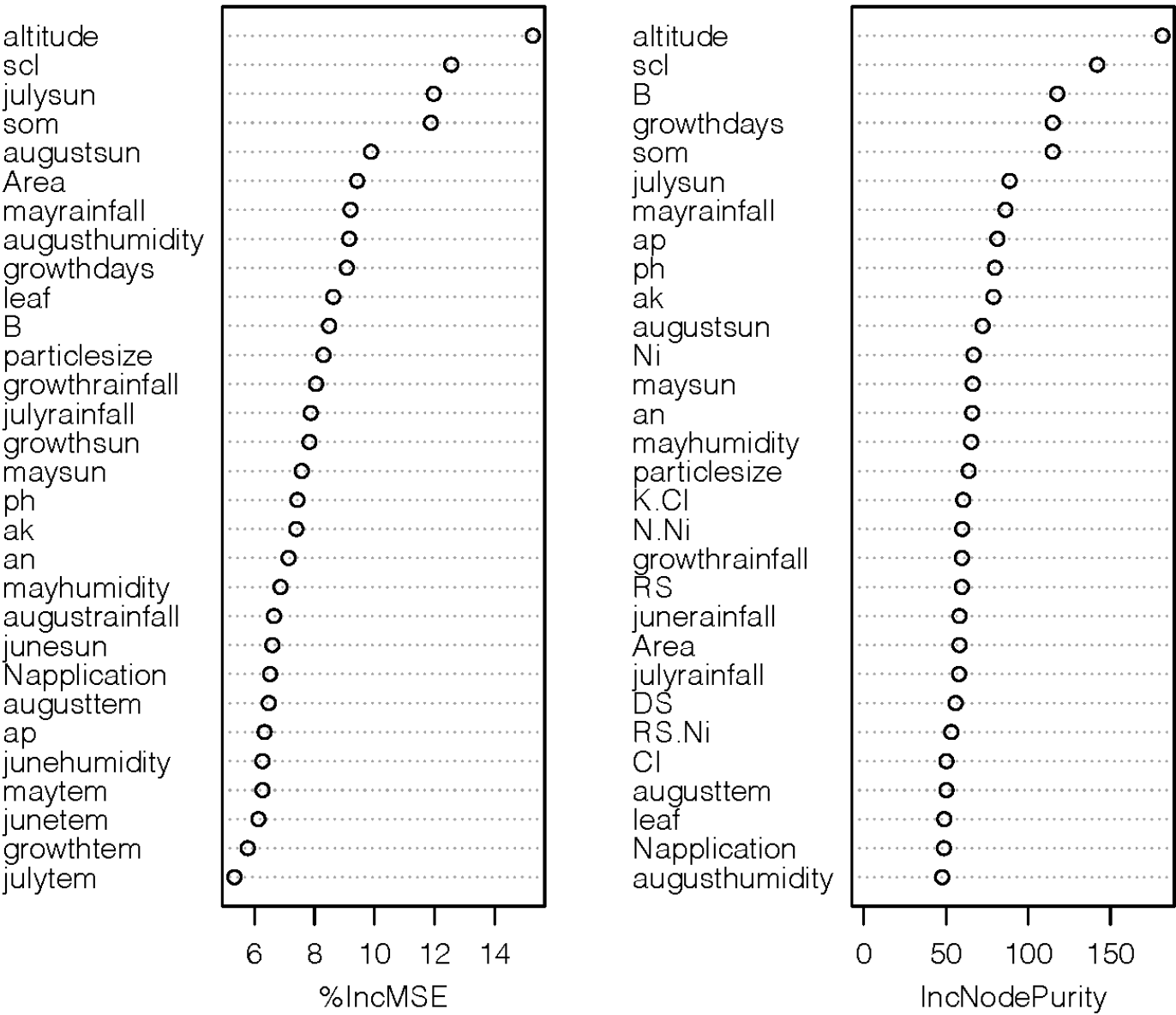
4.我想用主成分分析来分析数据，用哪些包？怎么用？

- 每种分析方法最好找对应的包，查询手册-[网络教程](#)-线上或线下寻求帮助
 - 需要从实践、学习中不断总结，大家分享
 - 好评包：dplyr、ggplot等
-

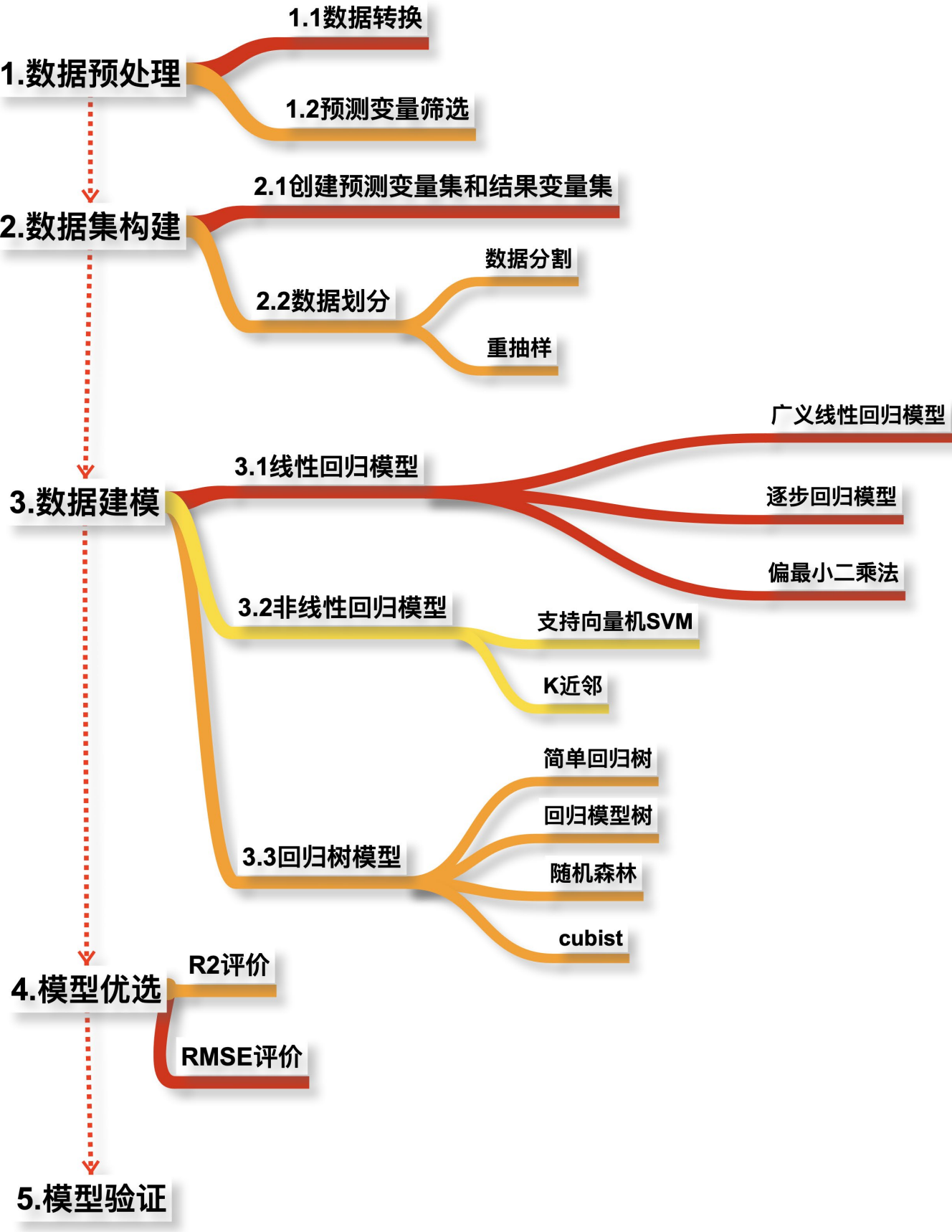
5.预测模型

随机森林、神经网络、支持向量机等，大量数据的训练，caret包

forest



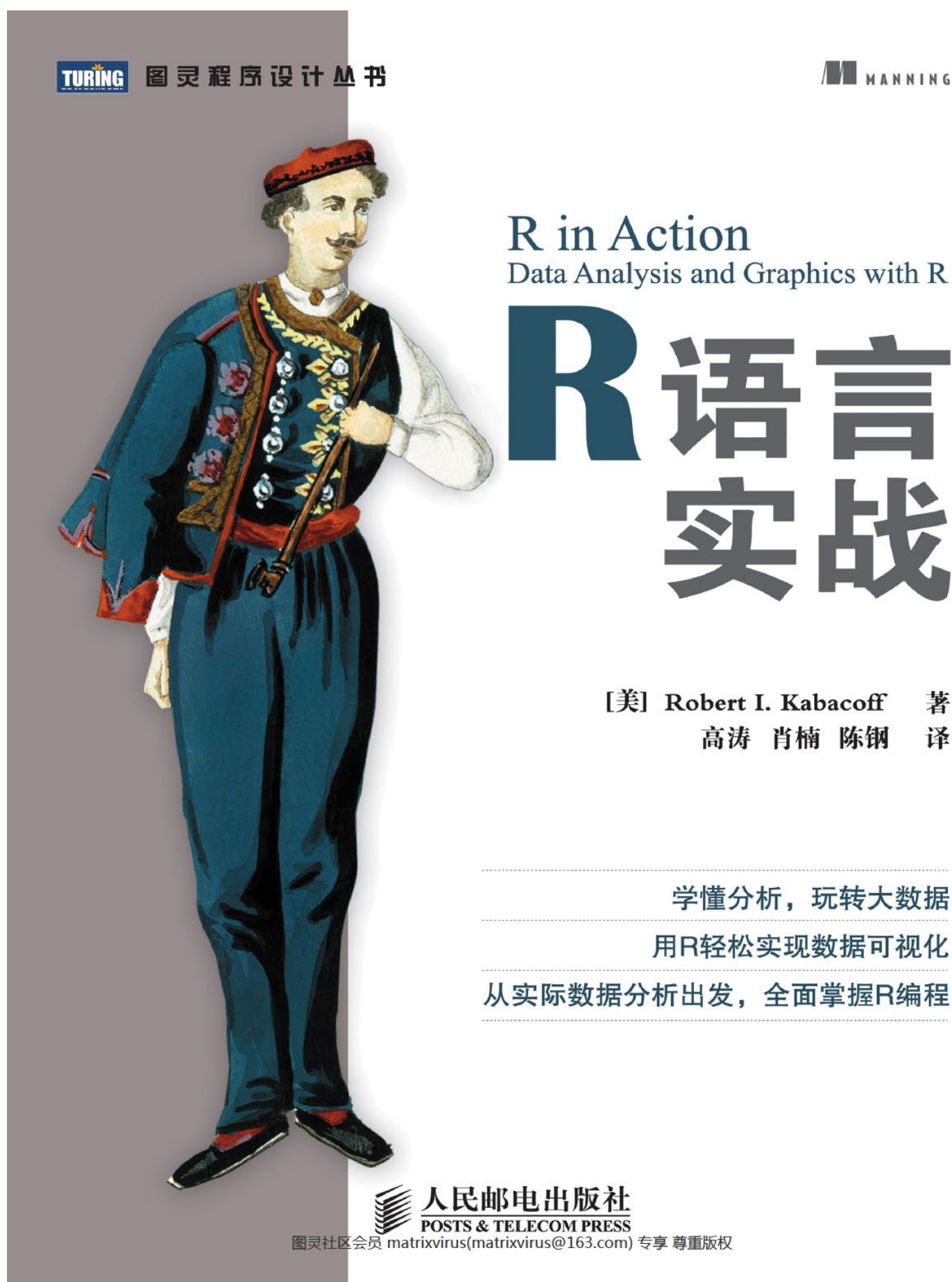
变量重要性排序



建模流程

6.Workflow

一切围绕《R语言实战》查询



R语言实战

6.1 软件安装

<https://www.r-project.org>

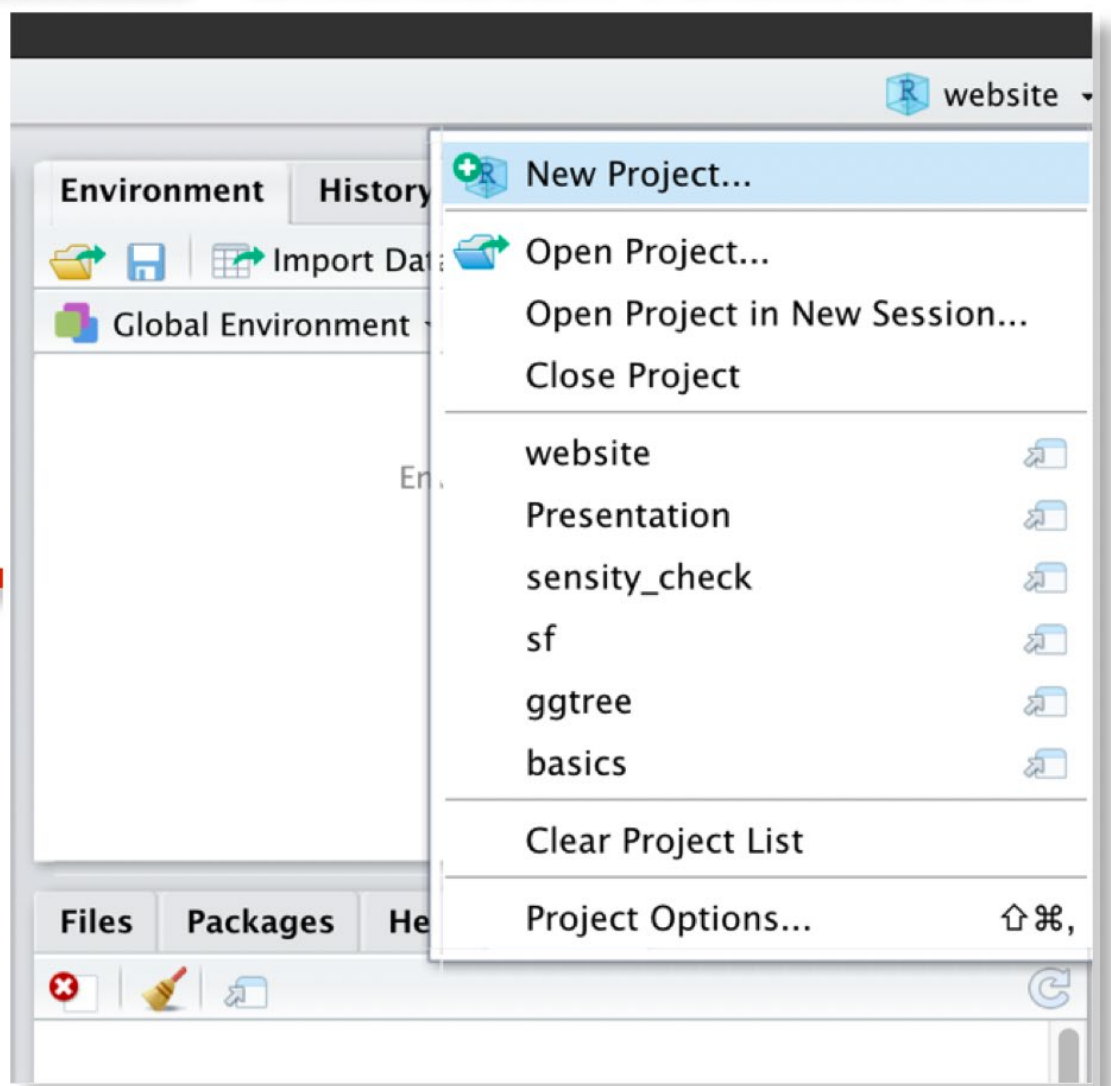
R安装

<https://www.rstudio.com/>

注意版本与R匹配

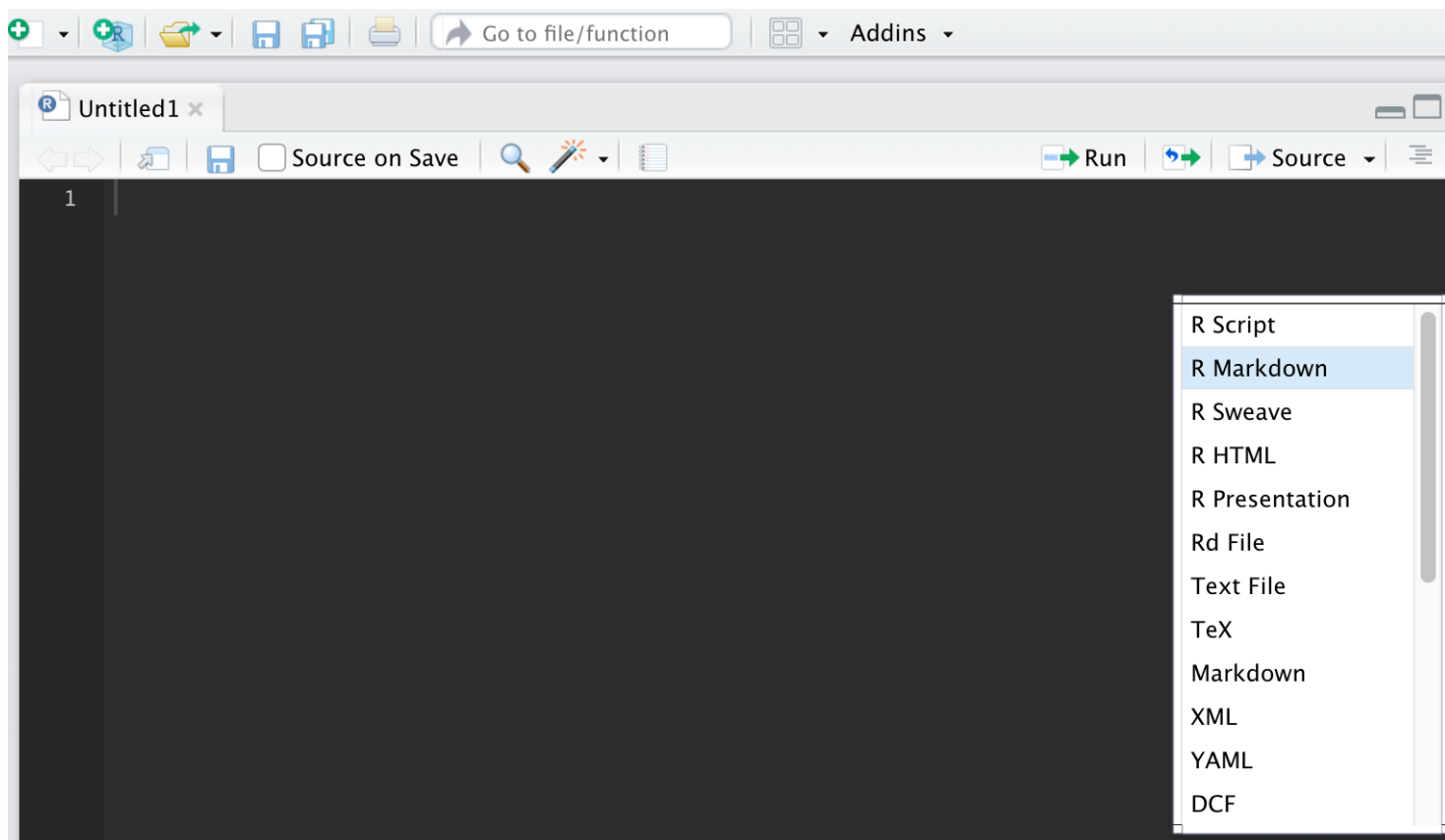
Rstudio安装

因为安装失败容易放弃，多在网上找答案，多尝试



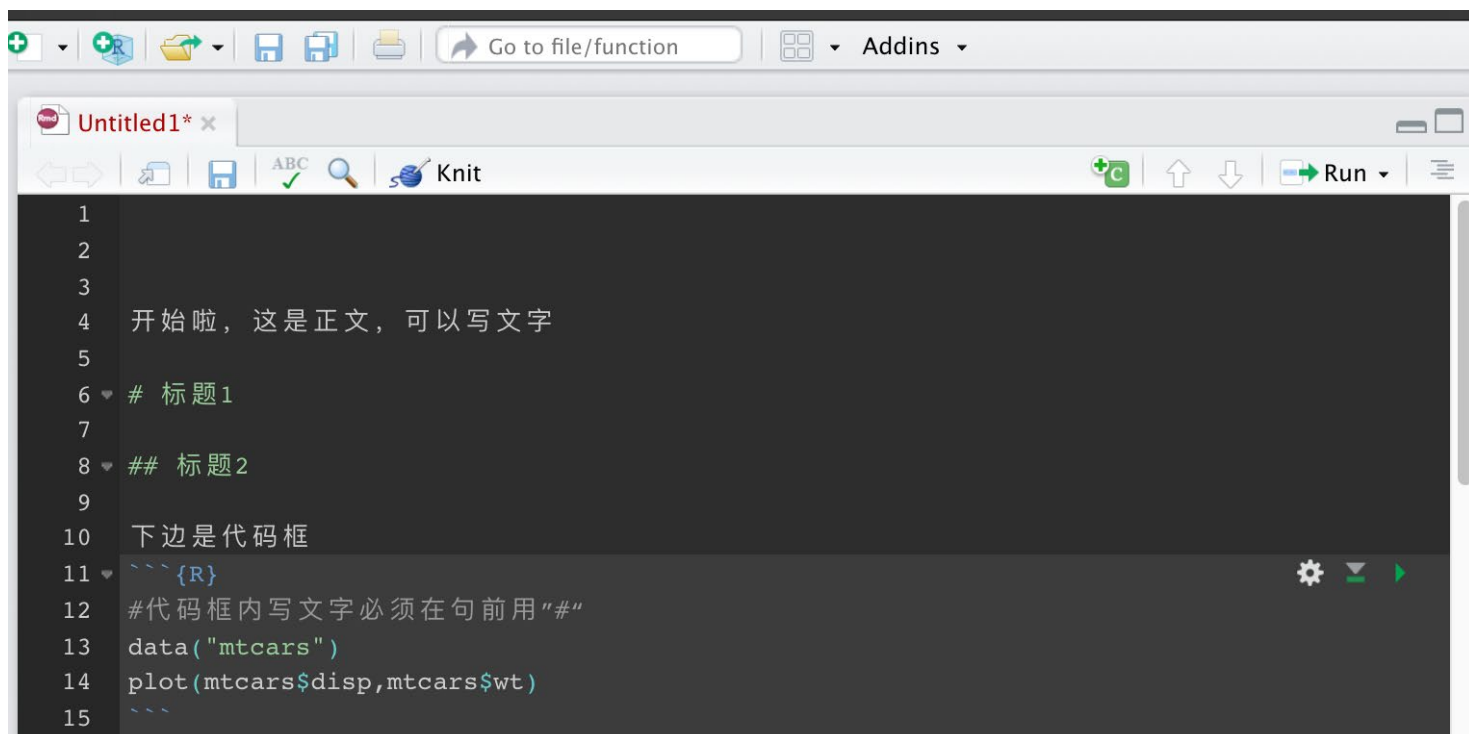
1

6.1 新建项目



2

6.2 使用RMarkdown



3

6.3 数据处理

6.3.1 导入数据

```
#install.packages("readxl")
library(readxl)
cardata <- read_excel("/Users/profits/Rdata/website/cardata.xlsx")
head(cardata)
```

```
## # A tibble: 6 x 12
##   Cartype      mpg   cyl  disp    hp  drat    wt   qsec    vs  am  gear  carb
##   <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Mazda RX4      21     6   160   110   3.9   2.62  16.5     0    1    4     4
## 2 Mazda RX4 W...  21     6   160   110   3.9   2.88  17.0     0    1    4     4
## 3 Datsun 710     22.8    4   108    93   3.85   2.32  18.6     1    1    4     1
## 4 Hornet 4 Dr...  21.4    6   258   110   3.08   3.22  19.4     1    0    3     1
## 5 Hornet Spor...  18.7    8   360   175   3.15   3.44  17.0     0    0    3     2
## 6 Valiant       18.1    6   225   105   2.76   3.46  20.2     1    0    3     1
```

6.3.3 筛选数据

筛选disp高于120，VS=1的数据

```
library(dplyr)
cardata_filter<-cardata%>%filter(dis>120,vs==1)
cardata_filter
```

```
## # A tibble: 8 x 12
##   Cartype      mpg   cyl  disp    hp  drat    wt   qsec    vs  am  gear  carb
##   <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Hornet 4 Dr...  21.4    6   258   110   3.08   3.22  19.4     1    0    3     1
## 2 Valiant       18.1    6   225   105   2.76   3.46  20.2     1    0    3     1
## 3 Merc 240D     24.4    4   147.    62   3.69   3.19  20      1    0    4     2
## 4 Merc 230      22.8    4   141.    95   3.92   3.15  22.9     1    0    4     2
## 5 Merc 280      19.2    6   168.   123   3.92   3.44  18.3     1    0    4     4
## 6 Merc 280C     17.8    6   168.   123   3.92   3.44  18.9     1    0    4     4
## 7 Toyota Coro...  21.5    4   120.    97   3.7    2.46  20.0     1    0    3     1
## 8 Volvo 142E    21.4    4   121   109   4.11   2.78  18.6     1    1    4     2
```

选取变量mpg和cyl

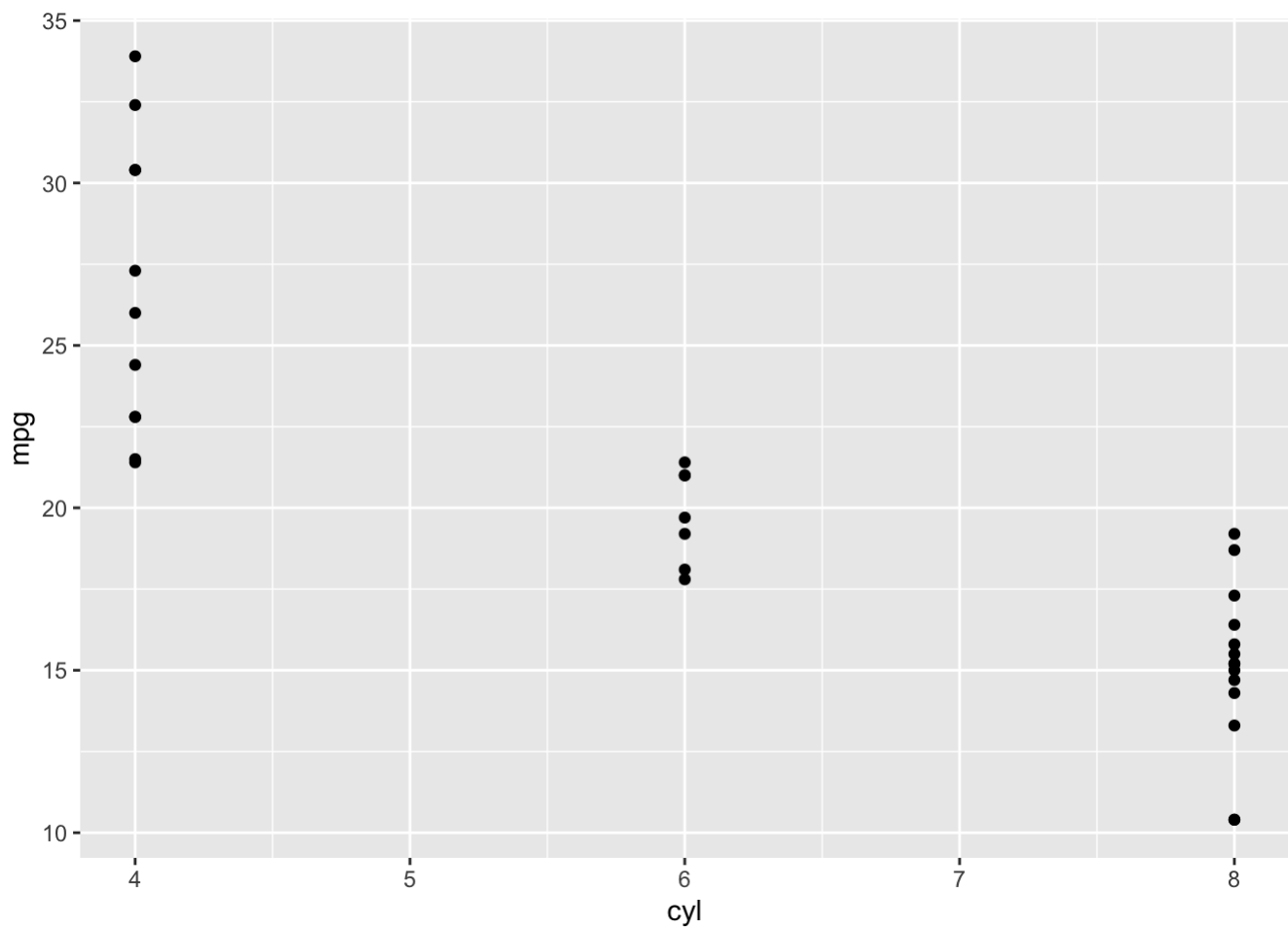
```
cardata_select<-cardata%>%select(mpg,cyl)
cardata_select
```

```
## # A tibble: 32 x 2
##      mpg    cyl
##    <dbl> <dbl>
##  1  21      6
##  2  21      6
##  3  22.8    4
##  4  21.4    6
##  5  18.7    8
##  6  18.1    6
##  7  14.3    8
##  8  24.4    4
##  9  22.8    4
## 10  19.2    6
## # ... with 22 more rows
```

6.3.4 数据展示

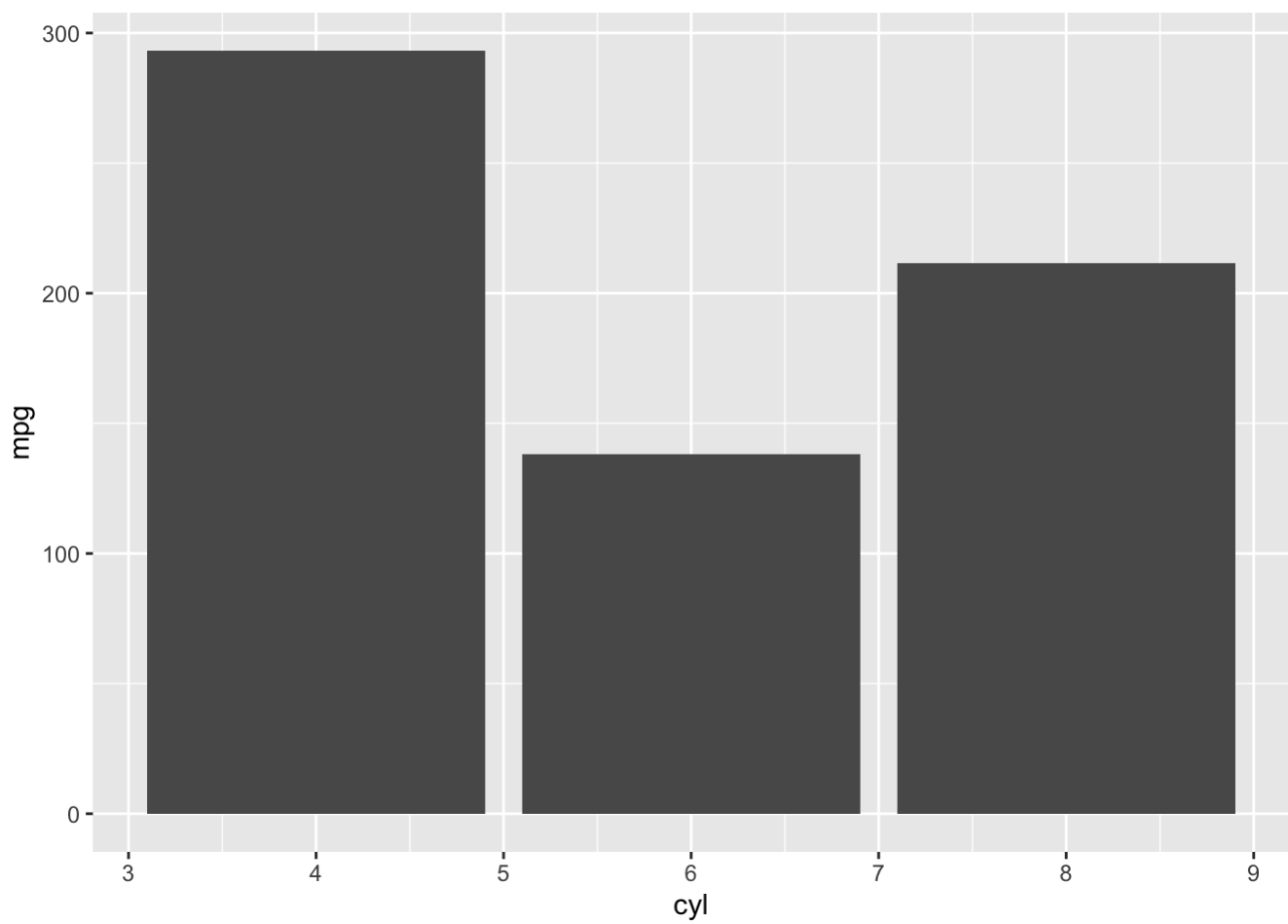
散点图

```
library(ggplot2)
ggplot(cardata,aes(x=cyl,y=mpg))+geom_point()
```



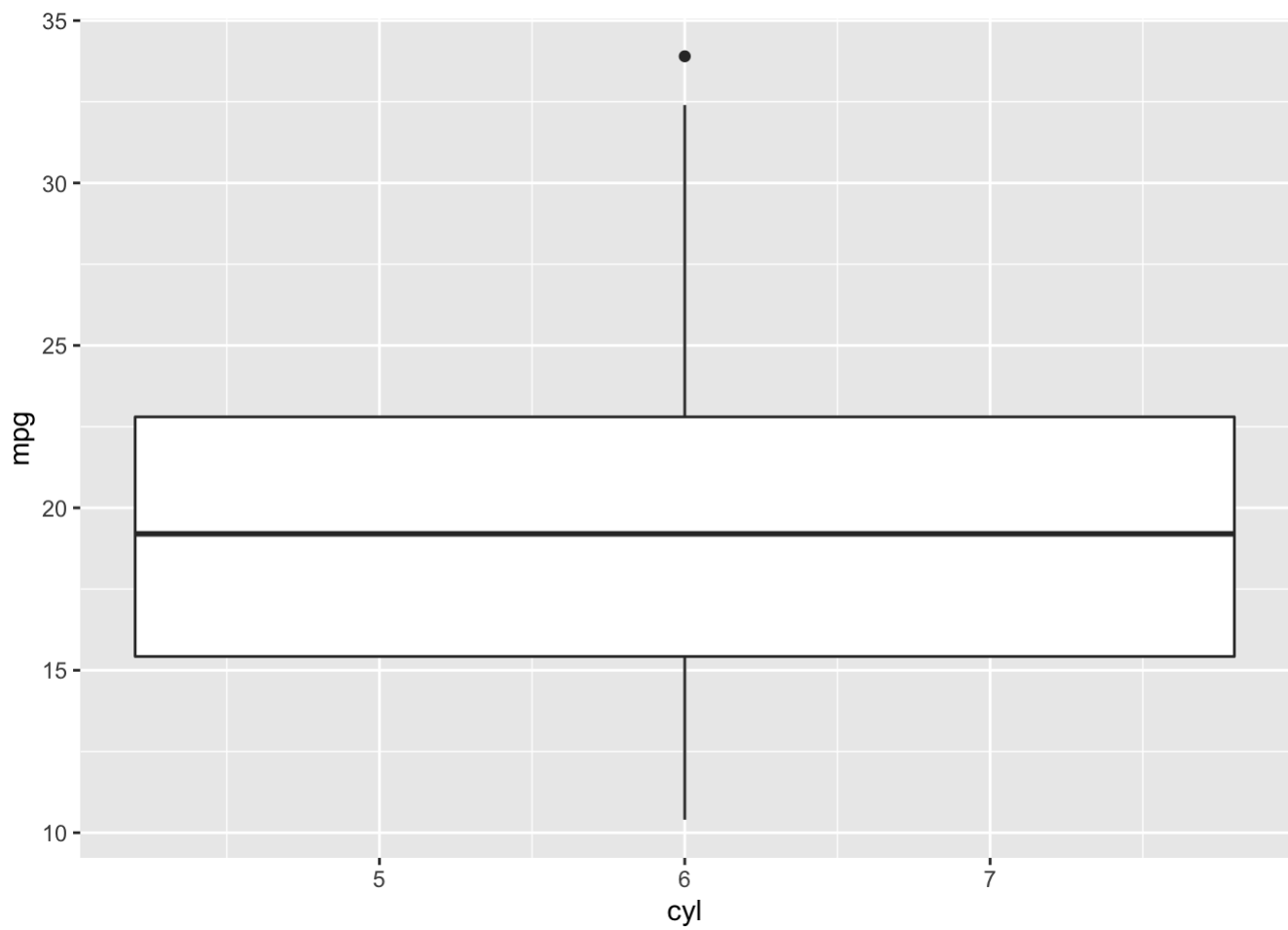
柱状图

```
ggplot(cardata,aes(x=cyl,y=mpg))+geom_col()
```



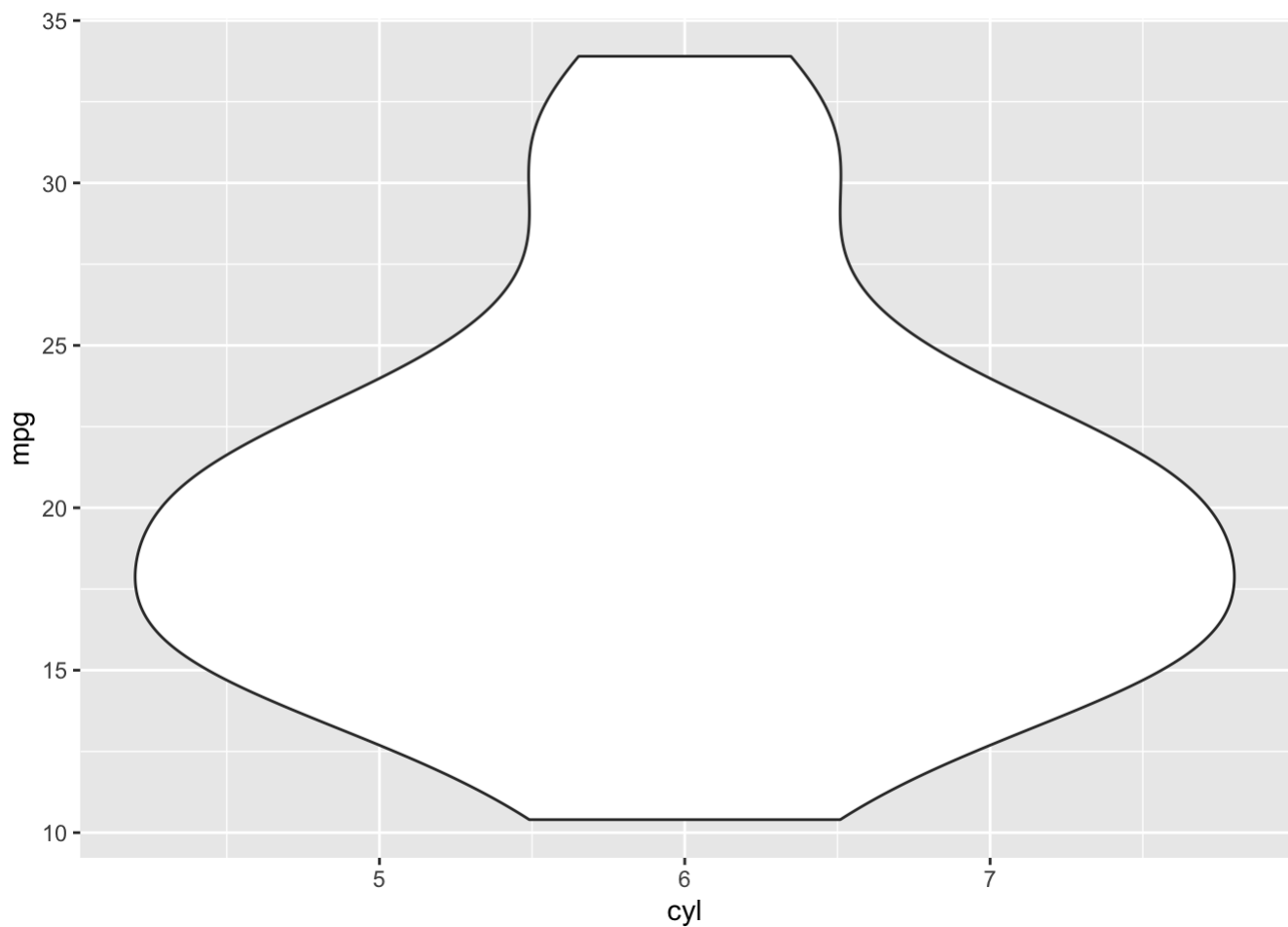
箱式图

```
ggplot(cardata,aes(x=cyl,y=mpg))+geom_boxplot()
```



小提琴图

```
ggplot(cardata,aes(x=cyl,y=mpg))+geom_violin()
```

6.3.5 数据分析

```
#anova
cardata$cyl<-as.factor(cardata$cyl)
anova_fit<-aov(cardata$mpg~cardata$cyl)
summary(anova_fit)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## cardata$cyl  2   824.8    412.4     39.7 4.98e-09 ***
## Residuals   29   301.3     10.4
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

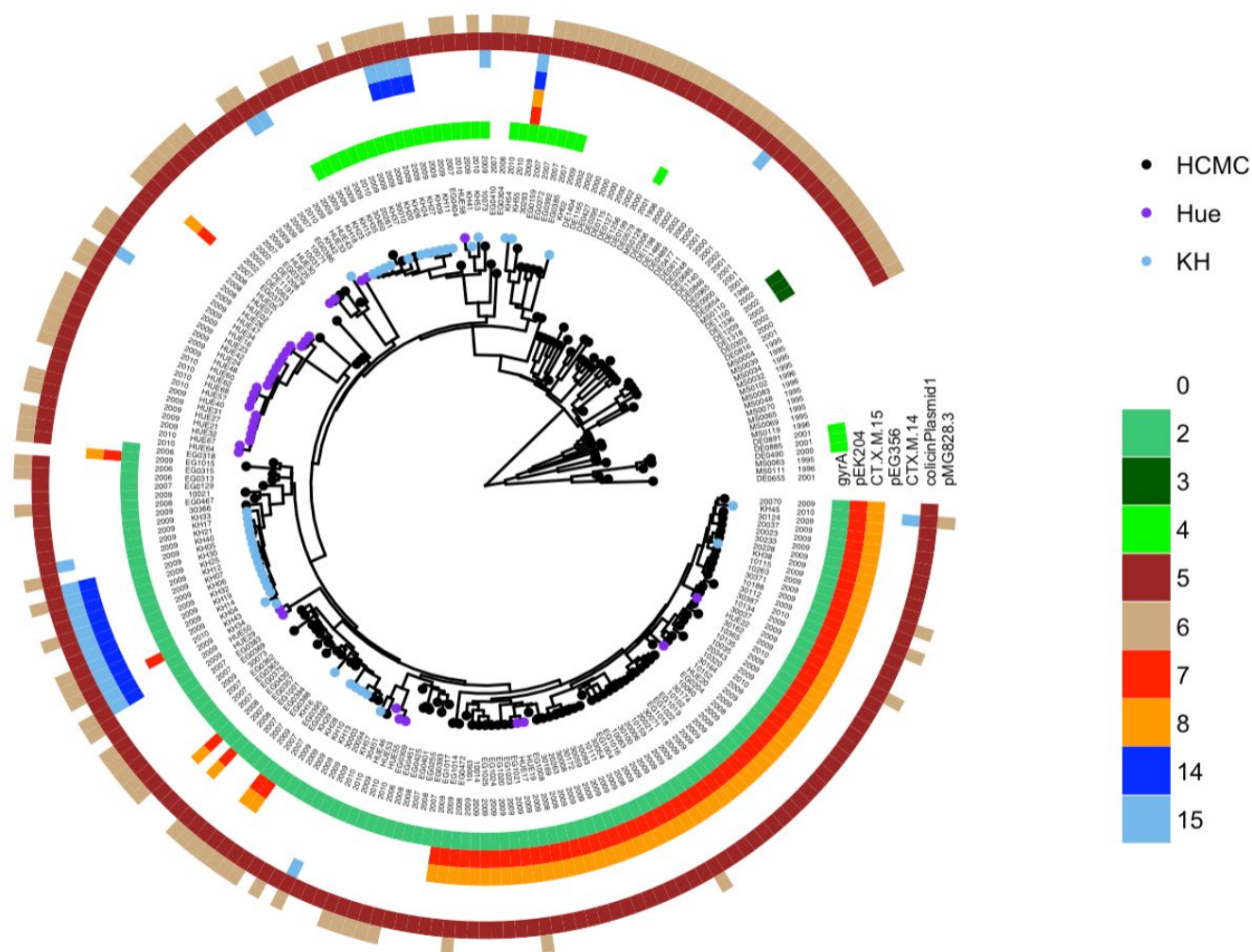
```
#lm_fit
lm_fit<-lm(mpg~cyl,cardata)
summary.aov(lm_fit)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
```

```
## cyl      2  824.8   412.4    39.7 4.98e-09 ***
## Residuals 29  301.3    10.4
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

6.4 更进一步的分析?

多重比较、线性混合效应模型、更好看的图，网上查找教程。比如，我想做个这样的聚类图：



网上找 [教程](#)

7.不断学习、分享、总结



只争朝夕-科研数据分析学习小组



该二维码 7 天内 (8月12日前) 有效, 重新进入将更新