





# 观数以形

艾新波 / 2018·北京



#### 课程体系









第3章 格言联璧话学习

🗐 第4章 源于数学、归于工程

中部: 执具

第5章 工欲善其事必先利其器

第6章 基础编程

第7章 数据对象









- 🗐 第11章 相随相伴、谓之关联

第12章 既是世间法、自当有分别

■ 第13章 方以类聚、物以群分

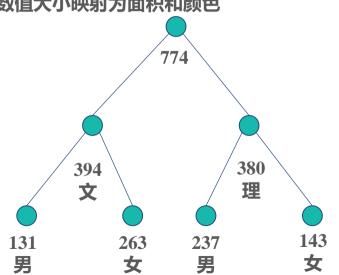
9 第14章 庐山烟雨浙江潮

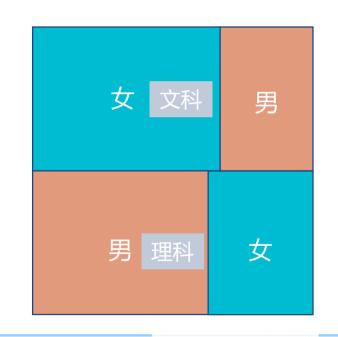
## 二维数据空间的形态

<b>类别</b>	刻画方法	图形示例	
类别变量vs类别变量 (因子vs因子)	矩形树图		0
数值变量vs数值变量 (向量vs向量)	散点图、 相关系数图		0
类别变量vs数值变量 (因子vs向量)	分组绘制箱线图/ 直方图/概率密度图等	FILE STATES OF THE PROPERTY OF	$\bigcirc$

#### 离散变量vs离散变量: 矩形树图

矩形树图采用嵌套的矩形来表达层级数据 数值大小映射为面积和颜色

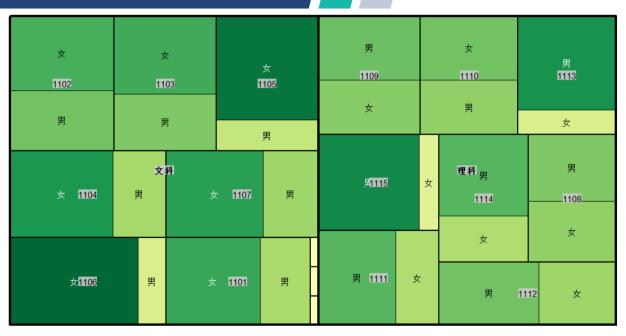




## 离散变量vs离散变量: 矩形树图

```
library(treemap)
cjb %>%
  group by (wlfk, bj, xb) %>%
  summarise(count = n()) %>%
  as.data.frame() %>%
  treemap(
    index = c("wlfk", "bj", "xb"),
    vSize = "count",
    vColor = "count",
    type = "value"
```

### 离散变量vs离散变量: 矩形树图



0 5 10 15 20 25 30 35 40 45

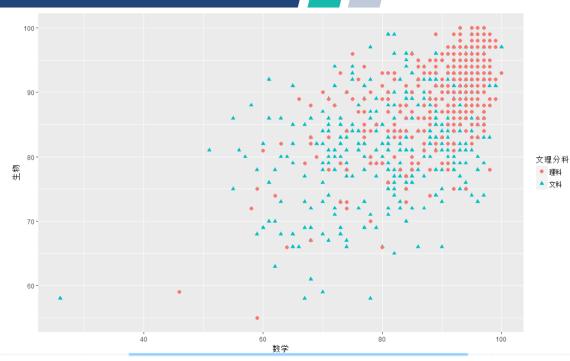
## 二维数据空间的形态

类别	刻画方法	图形示例	1 1 1 1 1
类别变量vs类别变量 (因子vs因子)	矩形树图		<b>②</b>
数值变量vs数值变量 (向量vs向量)	散点图、 相关系数图		С
类别变量vs数值变量 (因子vs向量)	分组绘制箱线图/ 直方图/概率密度图等	FIATOMIC TO THE CONTROL OF THE CONTR	С

#### 数值变量 vs 数值变量: 散点图

```
library(ggplot2)
ggplot(cjb,
      aes(x = sx,
           y = sw
           shape = wlfk,
           colour = wlfk)) +
  geom point(size = 2) +
  labs(x = "数学",
      v = "生物",
       colour = "文理分科",
       shape = "文理分科")
```

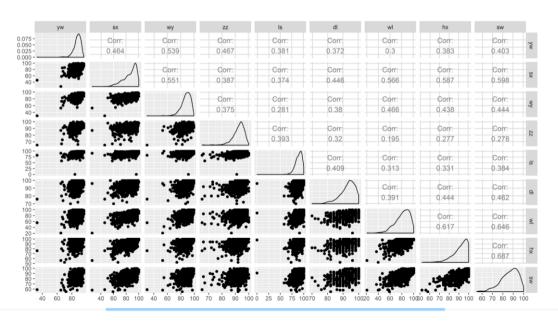
## 数值变量 vs 数值变量: 散点图



理科 ▲ 文科

#### 数值变量 vs 数值变量: 散点图矩阵

GGally::ggpairs(cjb, columns = 4:12)



#### 数值变量vs数值变量: 相关系数

#### 共变与相关:

当一个变量增大,另一个变量随之增大(或减小),称这种现象为共变,或相关(correlation)

两个变量之间具有共变现象, 称为具有相关关系 在没有特别说明的情况下, 相关系数一般是指线性相关系数

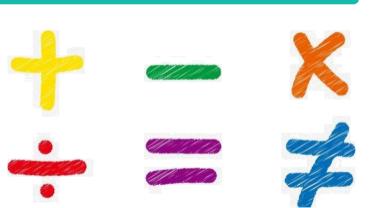
$$r_{xy} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{n} (y_i - \bar{y})^2}}$$

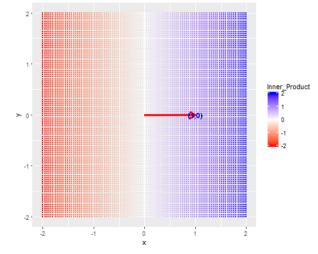
容易看到,相关系数的分子即为协方差(1/(n-1)) 顾名思义,协方差之"<mark>协"</mark>,在于刻画其变化的相似性

#### 数值变量vs数值变量: 相关系数

## 结合内积的概念,考虑协方差之"协"

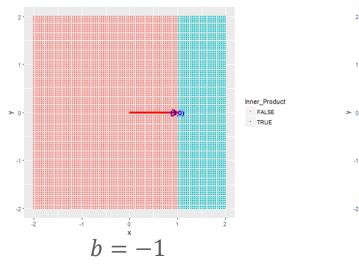
$$\langle x, y \rangle = x_1 y_1 + x_2 y_2 + \dots + x_n y_n$$

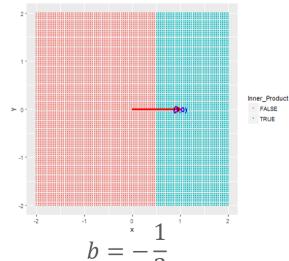




#### 数值变量vs数值变量: 相关系数

进一步延伸 $w^Tx + b = 0$ ,为下一步的分类超平面做准备





```
(cor coef <- cor(cjb[, 4:12]))
#>
     yw sx wy zz ls dl wl hx
                                             SW
#> yw 1.00 0.46 0.54 0.47 0.38 0.37 0.30 0.38 0.40
\#> sx 0.46 1.00 0.55 0.39 0.37 0.45 0.57 0.59 0.60
#> wy 0.54 0.55 1.00 0.37 0.28 0.38 0.47 0.44 0.44
\#> zz 0.47 0.39 0.37 1.00 0.39 0.32 0.20 0.28 0.28
#> 1s 0.38 0.37 0.28 0.39 1.00 0.41 0.31 0.33 0.38
#> d1 0.37 0.45 0.38 0.32 0.41 1.00 0.39 0.44 0.46
#> w1 0.30 0.57 0.47 0.20 0.31 0.39 1.00 0.62 0.65
#> hx 0.38 0.59 0.44 0.28 0.33 0.44 0.62 1.00 0.69
#> sw 0.40 0.60 0.44 0.28 0.38 0.46 0.65 0.69 1.00
```

```
(cor coef \leftarrow cor(cjb[, 4:12]))
cor coef %>%
  as.data.frame() %>%
  rownames to column (var = "km1") %>%
  gather(key = km2, value = cor num, -km1) %>%
  mutate(cor level = cut(cor num,
                   breaks= c(0, 0.3, 0.5, 0.8, 1),
                   right = FALSE)) %>%
  qqplot(aes(x = km1, y = km2, fill = cor level)) +
  geom tile(colour="white", size = 1.5) +
  geom text(aes(label = format(cor num, digits = 2))) +
  scale fill brewer(palette = "YlGn", name="相关系数区间")
```

```
km1 yw sx wy zz ls dl wl hx sw

1 yw 1.00 0.46 0.54 0.47 0.38 0.37 0.30 0.38 0.40

2 sx 0.46 1.00 0.55 0.39 0.37 0.45 0.57 0.59 0.60

3 wy 0.54 0.55 1.00 0.37 0.28 0.38 0.47 0.44 0.44

4 zz 0.47 0.39 0.37 1.00 0.39 0.32 0.20 0.28 0.28

5 ls 0.38 0.37 0.28 0.39 1.00 0.41 0.31 0.33 0.38

6 dl 0.37 0.45 0.38 0.32 0.41 1.00 0.39 0.44 0.46

7 wl 0.30 0.57 0.47 0.20 0.31 0.39 1.00 0.62 0.65

8 hx 0.38 0.59 0.44 0.28 0.33 0.44 0.62 1.00 0.69

9 sw 0.40 0.60 0.44 0.28 0.38 0.46 0.65 0.69 1.00
```

```
km1 km2 cor num cor level
               1.00
                          <NA>
    VW
        VW
               0.46 [0.3, 0.5)
    SX
        VW
               0.54 [0.5, 0.8)
    WV
        yw
    ΖZ
               0.47 [0.3, 0.5)
        VW
    ls
               0.38[0.3,0.5)
        VW
               0.37 [0.3, 0.5)
    d1
        VW
               0.30 [0.3, 0.5)
    wl
        VW
    hx
        VW
               0.38 [0.3, 0.5)
               0.40 [0.3, 0.5)
        VW
               0.46 [0.3, 0.5)
10
    VW
        SX
11
               1.00
                          <NA>
    SX
        SX
12
    WV
               0.55 [0.5, 0.8)
        SX
13
               0.39[0.3,0.5)
    ZZ
        SX
14
    1s
               0.37 [0.3, 0.5)
        SX
15
    d1
               0.45 [0.3, 0.5)
        SX
16
        SX
               0.57 [0.5, 0.8)
17
               0.59 [0.5, 0.8)
    hx
        SX
18
               0.60 [0.5, 0.8)
    SW
        SX
19
               0.54 [0.5, 0.8)
    yw
        Wy
20
               0.55 [0.5, 0.8)
        WV
```

```
(cor coef \leftarrow cor(cjb[, 4:12]))
cor coef %>%
  as.data.frame() %>%
  rownames to column (var = "km1") %>%
  gather(key = km2, value = cor num, -km1) %>%
  mutate(cor level = cut(cor num,
                   breaks= c(0, 0.3, 0.5, 0.8, 1),
                   right = FALSE)) %>%
  qqplot(aes(x = km1, y = km2, fill = cor level)) +
  geom tile(colour="white", size = 1.5) +
  geom text(aes(label = format(cor num, digits = 2))) +
  scale fill brewer(palette = "YlGn", name="相关系数区间")
```

zz -	0.32	0.28	0.39	0.28	0.39	0.20	0.37	0.47	1.00	
yw -	0.37	0.38	0.38	0.40	0.46	0.30	0.54	1.00	0.47	
wy -	0.38	0.44	0.28	0.44	0.55	0.47	1.00	0.54	0.37	
wl -	0.39	0.62	0.31	0.65	0.57	1.00	0.47	0.30	0.20	
ж sx -	0.45	0.59	0.37	0.60	1.00	0.57	0.55	0.46	0.39	
sw -	0.46	0.69	0.38	1.00	0.60	0.65	0.44	0.40	0.28	
ls -	0.41	0.33	1.00	0.38	0.37	0.31	0.28	0.38	0.39	
hx -	0.44	1.00	0.33	0.69	0.59	0.62	0.44	0.38	0.28	
dl -	1.00	0.44	0.41	0.46	0.45	0.39	0.38	0.37	0.32	
	dl	hx	ls	sw	sx km1	wl	wy	yw	zz	

相关系数区间

[0,0.3)

[0.3,0.5)

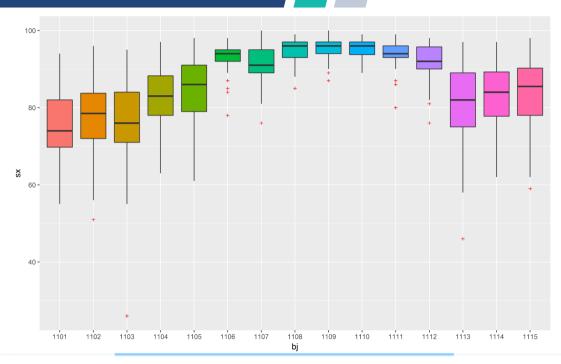
[0.5,0.8)

NA

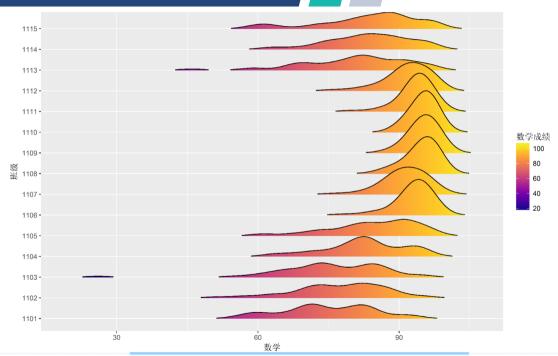
## 二维数据空间的形态

类别	刻画方法	图形示例
类别变量vs类别变量 (因子vs因子)	矩形树图	
数值变量vs数值变量 (向量vs向量)	散点图、 相关系数图	
类别变量vs数值变量 (因子vs向量)	分组绘制箱线图/ 直方图/概率密度图等	FIRST SERVICES TO CONTROL  TO

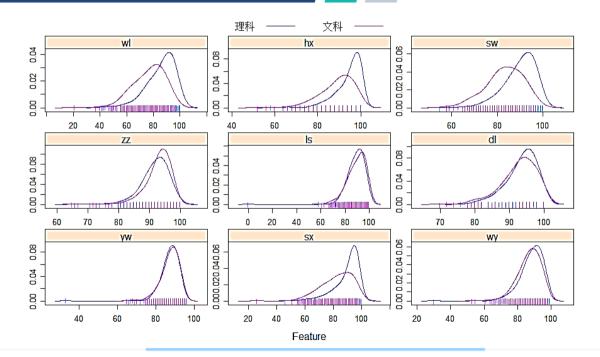
```
#分组绘制箱线图
#看看不同班级数学成绩的分布
library(ggplot2)
qqplot(cib, aes(x = bi,
               v = sx
               fill = bi)) +
 geom boxplot(outlier.colour = "red",
              outlier.shape = 3,
              outlier.size = 1) +
 labs(x = "班级", y = "数学成绩") +
  theme (legend.position = "none")
```



```
library(ggridges)#绘制层峦叠嶂图
library (viridis) #采用其中的颜色
qqplot(cjb, aes(x = sx, y = bj, fill = ..x..)) +
  geom density ridges gradient(
    scale = 2,
    rel min height = 0.01,
   qradient lwd = 1) +
  scale fill viridis(
   name = "数学成绩",
   option = "C") +
  labs(x = "数学", y = "班级")
```



```
#对于分类问题而言,在进行数据描述时
#最关键的,当属因变量vs自变量了
library (caret)
featurePlot(
 x = cib[, 4:12],
 y = cib$wlfk,
 plot = "density",
  scales = list(
   x = list(relation = "free"),
   y = list(relation = "free")),
 adjust = 1.5,
 pch = "|",
 auto.key = list(columns = 2))
```



#### 延伸阅读: CHEAT SHEET(ggplot2)

#### ONE VARIABLE continuous

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

## TWO VARIABLES continuous bivariate distribution

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

#### THREE VARIABLES

seals\$z <- 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

# 謝謝聆听 Thank you

#### 教师个人联系方式

艾新波

手机: 13641159546

QQ: 23127789

微信: 13641159546

E-mail: 13641159546@126.com

axb@bupt.edu.cn

地址:北京邮电大学科研楼917室

课程 网址: https://github.com/byaxb/RDataAnalytics



