

对象的5种基本类型

(1) 数值 (numeric:real numbers)

创建变量X，并将其赋值为1，赋值符号为左箭头加一个横杠

```
1 x<-1
2 x #查看变量x的值
3 [1] 1 #查看变量x的值返回的结果 [1]代表它后面接着的元素是x种的一个元素，没有方括号的1代表x中
4 class(x) #用class()函数查看变量x的类型,complex代表复数,character代表字符,integer代表整
5
6 #区分大小写
7 x<-5
8 X<-5
9
```

属性(attribute)

名称(name)

维度(dimensions: matrix,array)

类型(class)

长度(length)

向量(vector)

向量是一种可以包含多个元素的数据结构

(1) 使用vector()函数创建向量

```
1 vector(mode="logical",length= 0L)
2 #它包含两个参数：第一个参数是这个向量里参数的类型，第二个参数是这个向量包含的元素个数
3 x<-vector("character",length=10)
4 #变量x是字符型，有10个元素，每个元素是空的
```

(2) 直接在R中创建向量

```
1 x1<-1:3
2 # : 的意思是“到”，变量x1的内容是1, 2, 3
```

(3) 用c()函数创建向量

```
1 x2<-c(1,2,3)
2 #如果往c()函数中传入的元素类型不同的话，R会强制转换为同一类型的元素
3 x3<-c(TRUE,10,"a")#把三个不同类型的元素都转化为字符型元素
4 x4<-c("a","b","c")
5 #向量x4是一个字符类型向量，若想把它强制转换为数值型向量，需要用到as.numeric()函数
6 as.numeric(x4)
7 [1] NA NA NA #虽然强制把字符向量转化为数值型向量，但是R不知道如何把a,b,c转化为数字，因此用
```

一、矩阵 (matrix)

我们可以把矩阵看作：向量+维度属性，这个维度是一个整数向量，只包含两个元素行数和列数，分别用nrow和ncol表示

1.用matrix()函数创建矩阵

```
1 x<-matrix (nrow=3,ncol=2)
2 x<-matrix(1:6,nrow=3,ncol=2)#产生了一个3行2列，内容为1到6，且按列排列的矩阵
3 dim(x) #用矩阵维度的属性dim()查看矩阵的行列数
4 #返回 [1] 3 2
5 attributes(x) #查看矩阵x有多少属性，有哪些属性
```

2.用向量+维度属性创建矩阵

```
1 #创建一个向量
2 y<-1:6
3 #给向量添加维度信息
4 dim(y)<-c(2,3) #创建出一个2行3列的矩阵y
```

3.拼接矩阵

```
1 #创建一个矩阵
2 y2<-matrix(1:6,nrow=2,ncol=3)
3
4 #用rbind()函数按行拼接矩阵
5 rbind(y,y2)
6 #由于矩阵y和y2都是2行3列的矩阵，那么按行拼接就是4行3列的矩阵
7
8 #用cbind()函数按列拼接矩阵
9 cbind(y,y2)
10 #由于矩阵y和y2都是2行3列的矩阵，那么按列拼接就是2行6列的矩阵
```

二、数组 (array)

数组与矩阵类似，但是数组的维度可以大于2

1、用array()函数创建2维数组

```
1 > x<-array(1:24,dim=c(4,6))
2 > x
3      [,1] [,2] [,3] [,4] [,5] [,6]
4 [1,]    1    5    9   13   17   21
5 [2,]    2    6   10   14   18   22
6 [3,]    3    7   11   15   19   23
7 [4,]    4    8   12   16   20   24
```

2、用array()函数创建3维数组

```
1 > x1<-array(1:24,dim=c(2,3,4))
2 > x1
3 , , 1
4
5      [,1] [,2] [,3]
6 [1,]    1    3    5
7 [2,]    2    4    6
8
9 , , 2
10
11     [,1] [,2] [,3]
12 [1,]    7    9   11
13 [2,]    8   10   12
14
15 , , 3
16
17     [,1] [,2] [,3]
18 [1,]   13   15   17
19 [2,]   14   16   18
20
21 , , 4
22
23     [,1] [,2] [,3]
24 [1,]   19   21   23
25 [2,]   20   22   24
```


一、创建列表

向量、矩阵和数组，它们的共性是可以包含多个元素，并且元素和元素之间的类型必须是相同的，列表和它们最大的区别是：它可以包含不同类型的对象

1.用list()函数创建列表

```
1 > m<-list("a","1","2L",1+2i,FALSE)
2 > m
3 [[1]]
4 [1] "a"
5
6 [[2]]
7 [1] "1"
8
9 [[3]]
10 [1] "2L"
11
12 [[4]]
13 [1] 1+2i
14
15 [[5]]
16 [1] FALSE
17
18 #该列表共包含5个元素，它们分别是字符型、数值型、整数型、复数型和逻辑型这5种类型的对象
```

2.用list()函数给列表里的元素命名

```
1 > a<-list(a=1,b=2,c=3)
2 > a
3 $a
4 [1] 1
5
6 $b
7 [1] 2
8
9 $c
10 [1] 3
```

```
11
12 #该列表包含3个元素a,b,c,它们内容分别是1, 2, 3
```

3.用list()函数创建列表中的每一个元素包含的元素个数大于1的列表

```
1 > c<-list(c(1,2,3),c(4,5,6,7,8))
2 > c
3 [[1]]
4 [1] 1 2 3
5
6 [[2]]
7 [1] 4 5 6 7 8
```

4.引入矩阵的维度

(1) 先用matrix()函数创建一个矩阵

```
1 > x<-matrix(1:6,nrow=2,ncol=3) #1: 6先按列排列再按照行
2 > x
3      [,1] [,2] [,3]
4 [1,]    1    3    5
5 [2,]    2    4    6
```

(2) 再给矩阵的行列命名

```
1 > dimnames(x)<-list(c("a","b"),c("c","d","e"))
2 > x
3      c d e
4 a 1 3 5
5 b 2 4 6
```

二、构建列表的子集

1、先用list()函数构建一个列表

```
1 > x<-list(id=1:3,height=180,gender="male")
2 > x
3 $id
4 [1] 1 2 3
5
6 $height
7 [1] 180
8
9 $gender
10 [1] "male"
```

2、查看列表第一个元素的名称及内容

```
1 #两种方法
2 > x[1]
3 $id
4 [1] 1 2 3
5
6 > x["id"]
7 $id
8 [1] 1 2 3
```

3、只查看列表第一个元素的内容

```
1 #三种方法
2 > x[[1]]
3 [1] 1 2 3
4 > x[["id"]]
5 [1] 1 2 3
6 > x$id
7 [1] 1 2 3
```

4、查看列表中多个元素的名称及内容

```
1 > x[c(1,3)]
```



```

2 $id
3 [1] 1 2 3
4 $gender
5 [1] "male"
6
7 > x[c(1,2)]
8 $id
9 [1] 1 2 3
10 $height
11 [1] 180
12
13 > x[c(1,2,3)]
14 $id
15 [1] 1 2 3
16 $height
17 [1] 180
18 $gender
19 [1] "male"

```

一个很容易犯错的小知识点

```

1 > y<-"id"
2 > x[["id"]]
3 [1] 1 2 3
4 > x[[y]]
5 [1] 1 2 3
6 > x$id
7 [1] 1 2 3
8 > x$y
9 NULL
10 #使用嵌套的方括号[]能够引用包含了名称的变量y
11 #x$y返回的是空值，所以我们使用$符号只能直接引用名称"id"，不能引用包含了名称的变量y

```

5、从列表中获取嵌套的元素

(1)先创建一个列表

```

1 > x<-list(a=list(1,2,3),b=c("Monday","Tuesday"))
2 > x

```

```

3 $a
4 $a[[1]]
5 [1] 1
6
7 $a[[2]]
8 [1] 2
9
10 $a[[3]]
11 [1] 3
12
13
14 $b
15 [1] "Monday" "Tuesday"

```

(2) 提取列表中的列表的元素的内容

使用两个嵌套的方括号[]提取

```

1 > x[[1]][[1]]
2 [1] 1
3 > x[[1]][[2]]
4 [1] 2
5 > x[[1]][[3]]
6 [1] 3

```

```

1 > x[[c(1,1)]]
2 [1] 1
3 > x[[c(1,2)]]
4 [1] 2
5 > x[[c(1,3)]]
6 [1] 3
7 > x[[c(2,1)]]
8 [1] "Monday"
9 > x[[c(2,2)]]
10 [1] "Tuesday"
11
12 #c()函数中的第一个参数2表示列表x中的第二个元素，第二个参数1表示第二个元素中的第1个元素

```


采用

一、判断缺失值

缺失值有两种表示方式：NA和NaN，两者之间的关系是：NaN属于NA，NA不属于NaN。原因是：NaN一般表示数值类型的缺失值，NA表示的数据类型可以多种，比如：数值的缺失值、字符的缺失值等。

1.先用c()函数创建一个向量

```
1 > x<-c(1,NA,NA,2,3)
2 > x
3 [1] 1 NA NA 2 3
```

2.用is.na()和is.nan()函数判断缺失值

```
1 > is.na(x)
2 [1] FALSE TRUE TRUE FALSE FALSE
3 > is.nan(x)
4 [1] FALSE FALSE FALSE FALSE FALSE
```

从结果可以看出：is.na()函数判断出了每个元素是否为缺失值，缺失值返回了真，数值返回了假；is.nan()函数没有判断出缺失值，说明在向量x中不存在NaN这种类型的缺失值，同时也印证了NA不属于NaN。若我们把向量x里的NA改成NaN，结果如下：

```
1 > x<-c(1,NaN,NaN,2,3)
2 > x
3 [1] 1 NaN NaN 2 3
4 > is.na(x)
5 [1] FALSE TRUE TRUE FALSE FALSE
6 > is.nan(x)
7 [1] FALSE TRUE TRUE FALSE FALSE
```

二、处理缺失值

```

1 #1.先用c()函数创建一个包含缺失值的向量
2 > x<-c(1,NA,NA,2,3)
3 > x
4 [1] 1 NA NA 2 3
5 #2.取向量中不是缺失值的元素，！的意思是取反
6 > x[!is.na(x)]
7 [1] 1 2 3
8
9 #3.用complete.case()函数取多个向量对应位置都不是缺失值的元素
10 > x<-c(1,NA,NA,2,3)
11 > y<-c(NA,"a","b","c","d")
12 > z<-c(1L,2L,3L,NA,4L)
13 > w<-complete.cases(x,y,z)
14 > w
15 [1] FALSE FALSE FALSE FALSE TRUE
16 #返回结果是逻辑向量，其中只有x和y对应位置都不是缺失值的才会返回TRUE，否则返回FALSE。我们用x[

```

```

1 > x[w]
2 [1] 3
3 > y[w]
4 [1] "d"
5 > z[w]
6 [1] 4

```

三、实例

1.加载数据集所在的包

```

1 > library(datasets)

```

2.用head()函数查看数据集的前6行

```

1 head(airquality)
2   Ozone Solar.R Wind Temp Month Day
3 1    41    190  7.4  67     5    1

```

4	2	36	118	8.0	72	5	2
5	3	12	149	12.6	74	5	3
6	4	18	313	11.5	62	5	4
7	5	NA	NA	14.3	56	5	5
8	6	28	NA	14.9	66	5	6

3.用complete.cases()函数选择在变量

```

1 g<-complete.cases(airquality)
2 > g
3 [1] TRUE TRUE TRUE TRUE FALSE FALSE TRUE TRUE TRUE FALSE
4 [11] FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
5 [21] TRUE TRUE TRUE TRUE FALSE FALSE FALSE TRUE TRUE TRUE
6 [31] TRUE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE
7 [41] TRUE FALSE FALSE TRUE FALSE FALSE TRUE TRUE TRUE TRUE
8 [51] TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
9 [61] FALSE TRUE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE
10 [71] TRUE FALSE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE
11 [81] TRUE TRUE FALSE FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
12 [91] TRUE TRUE TRUE TRUE TRUE FALSE FALSE FALSE TRUE TRUE
13 [101] TRUE FALSE FALSE TRUE TRUE TRUE FALSE TRUE TRUE TRUE
14 [111] TRUE TRUE TRUE TRUE FALSE TRUE TRUE TRUE TRUE FALSE TRUE
15 [121] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
16 [131] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
17 [141] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE FALSE
18 [151] TRUE TRUE TRUE
19

```

4.选择该数据集不包含缺失值记录的前10行，且每个变量都要

```

1 > airquality[g,][1:10,]
2   Ozone Solar.R Wind Temp Month Day
3 1    41    190  7.4  67    5    1
4 2    36    118  8.0  72    5    2
5 3    12    149 12.6  74    5    3
6 4    18    313 11.5  62    5    4
7 7    23    299  8.6  65    5    7
8 8    19     99 13.8  59    5    8

```

9	9	8	19	20.1	61	5	9
10	12	16	256	9.7	69	5	12
11	13	11	290	9.2	66	5	13
12	14	14	274	10.9	68	5	14

四、识别缺失值的模式

```
1 install.packages("mice")
2 library(mice)
```

```
1 #用md.pattern()函数查看数据缺失值的分布
2 md.pattern(data)
```

五、处理缺失值

1.行删除法

用na.omit()函数删除不完整观测

```
1 library(datasets)
2 > data<-head(airquality)
3 > data
4   Ozone Solar.R Wind Temp Month Day
5 1    41    190  7.4  67     5   1
6 2    36    118  8.0  72     5   2
7 3    12    149 12.6  74     5   3
8 4    18    313 11.5  62     5   4
9 5    NA     NA 14.3  56     5   5
10 6    28     NA 14.9  66     5   6
11 > newdata<-na.omit(data) #用na.omit()函数删除不完整观测
12 > newdata
13   Ozone Solar.R Wind Temp Month Day
14 1    41    190  7.4  67     5   1
15 2    36    118  8.0  72     5   2
16 3    12    149 12.6  74     5   3
17 4    18    313 11.5  62     5   4
```

2.多重插补法

多重插补法（MI）是一种基于重复模拟的处理缺失值的方法，多重插补是从一个包含缺失值的数据集中生成一组完整的数据集。每个模拟数据集中，缺失数据将使用蒙特卡洛方法来填补。

```
1 library(datasets)
2 data<-head(airquality)
3 data
4
5 data1<-mice(data,m=6)
6 #用mice()函数从包含缺失数据的数据框开始，进行6重插补，即生成6个完整数据集
7 #mice()函数的第一个参数data为数据，第二个参数m为要返回的完整数据集的个数
8
9 fit<-with(data1,lm(Solar.R~Wind+Temp+Month+Day+Ozone))
10 fit
11 #用with()函数依次对6个完整数据集分别应用lm()模型
12 #结果分别返回6个完整数据集的回归结果
13
14 jiegua<-pool(fit)
15 jiegua
16 summary(jiegua)
17 #用pool()函数汇总回归结果
18
19 result<-complete(data1,action=3)
20 result
21 #选择第三个插补数据集作为结果
```

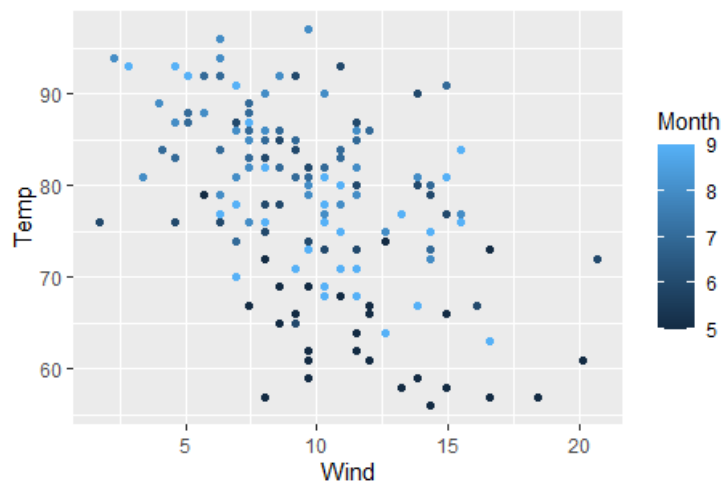
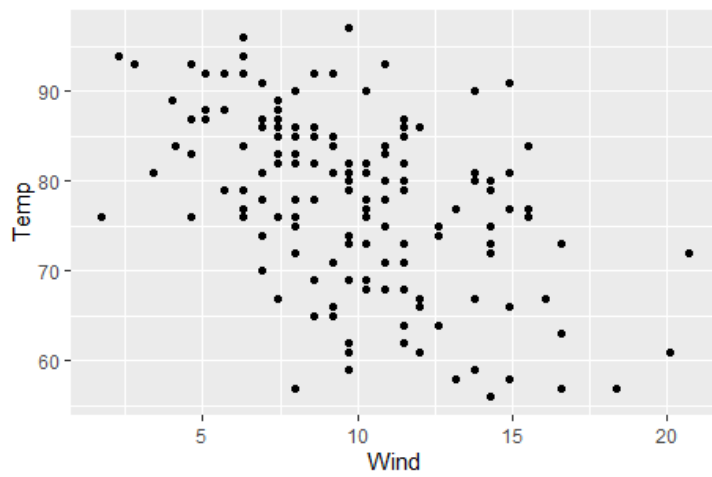


```
1 a<-read.csv()
2 #更改工作空间目录，直接用文件名读入数据
3 #（1）用read.csv()读入数据
4 #（2）用setwd()函数重新设定工作空间
5 > setwd("D:\\")
6 #（3）用getwd()函数查看目前工作空间所在位置
7 > getwd()
8 #（4）再次用read.csv()函数读入数据
9 > a<-read.csv()
10
11
12
13
14
15 #安装readxl程序包
16 install.packages("readxl")
17 library(readxl)
18 #用read_excel()函数读入数据
19 data=read_excel()
20
21 #表格状数据的读入，用read.table()函数
22 zhushi<-read.table("竹石.txt")
23 zhushi<-read.table("竹石.txt",header=T) #数据中第一行的数据作为列数据的变量名
24
25 #读取没有任何规则的数据需要readlines()函数
26 shuju1=readlines("数据小说数据.txt")
```

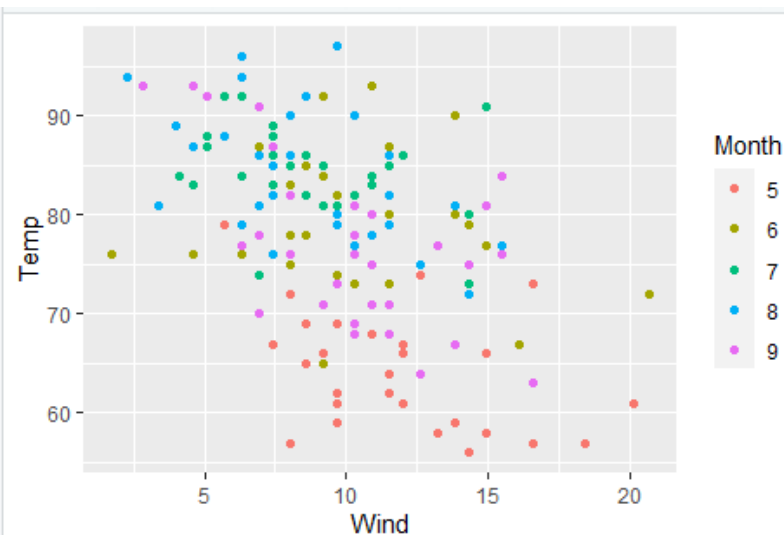
层 (layer)

Data	感兴趣的变量 (data frame)
Aesthetics(美学层)	x-axis(x轴) /y-axis(y轴) /color(颜色) /fill(填充的颜色) /size(大小) /labels(标签) /alpha(透明度) /shape(形状) /linearwidth(线宽) /lineartype(线的类型)
Geometries(几何客体层)	point(散点图) /line(线图) /histogram(柱状图) /bar(条形图) /boxplot(箱图)
Facets (划分绘图面板)	columns(行) ,rows(行)
Statistics (统计层) (目的: 添加统计信息)	binning/smoothing/descriptive/inferential
coordinates(坐标系)	cartesian/fixed/polar/limits
Themes(主题)	non-data ink(和数据无关的风格设计)

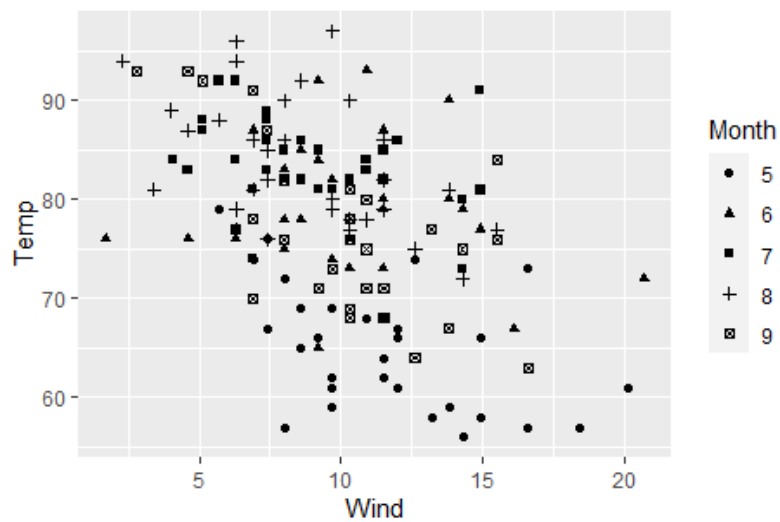
```
1 install.packages("ggplot2")
2 library(ggplot2)
3 #绘制风速~温度的散点图
4 qplot(Wind,Temp,data=airquality)
5 #按月份显示不同的颜色
6 qplot(Wind,Temp,data=airquality,color=Month)
```



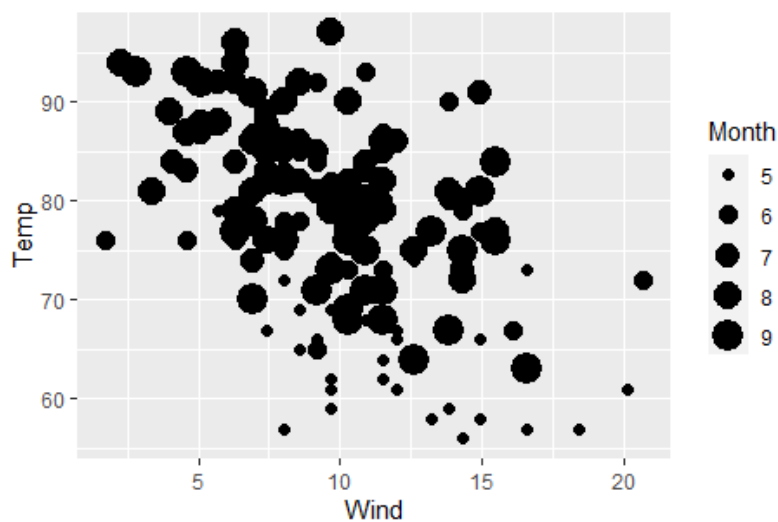
```
1 #由于元数据集中Month变量是一个连续变量，因此颜色为渐变色。如果想要不同月份显示不同颜色（非连续的渐变色），需要
2 airquality$Month<-factor(airquality$Month)
3 qplot(Wind,Temp,data=airquality,color=Month)
```



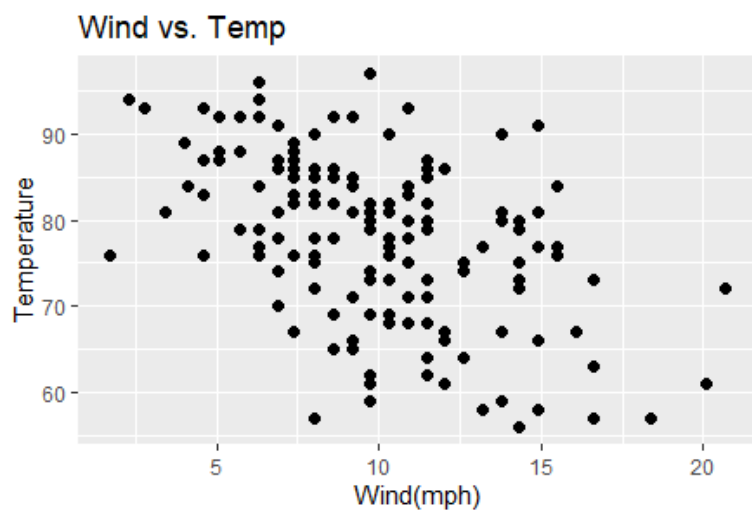
```
1 #不同月份显示不同形状
2 qplot(Wind,Temp,data=airquality,shape=Month)
```



```
1 #不同月份显示不同的大小
2 qplot(Wind,Temp,data=airquality,size=Month)
```



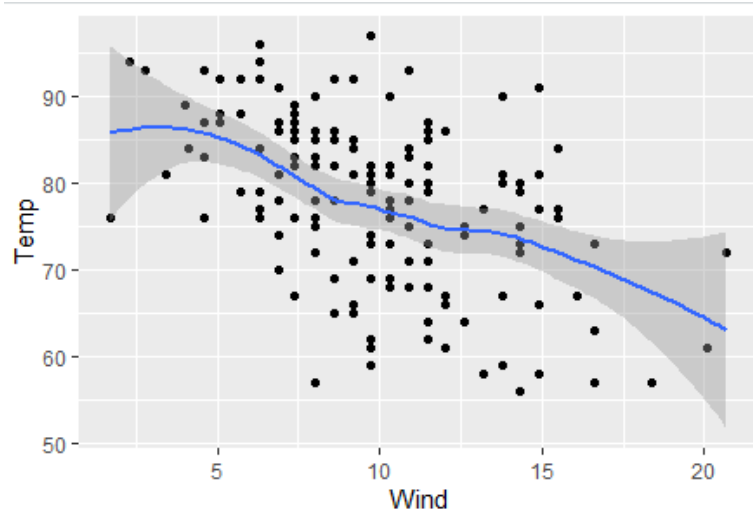
```
1 #设置大小为2, x,y轴标签分别为“Wind(mph)”,“emperature”,标题为“Wind vs. Temp”
2 qplot(Wind,Temp,data=airquality,size=I(2),main="Wind vs. Temp",xlab="Wind(mph)",ylab="Temperature")
```



```

1 #生成带平滑回归线的散点图
2 qplot(Wind,Temp,data=airquality,geom = c("point","smooth"))
3 #代码中的smooth会给一条根据point(点)拟合出来的回归线，它可以算作统计信息，
4 #可以看出，图中多了一条按照默认方法拟合出的平滑的蓝色曲线(回归线)和灰色的条块（置信区间）

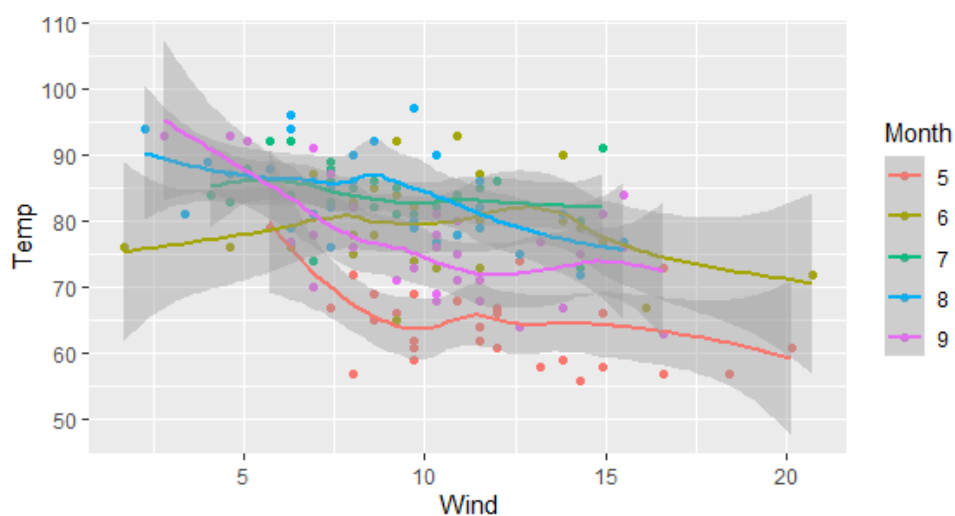
```



```

1 #按月份生成不同颜色的平滑回归线
2 qplot(Wind,Temp,data=airquality,geom = c("point","smooth"),color=Month)

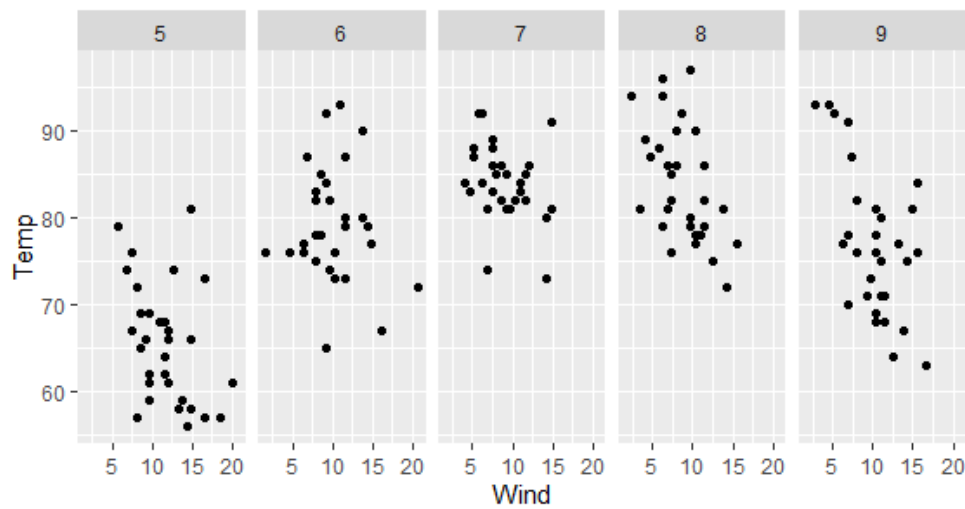
```



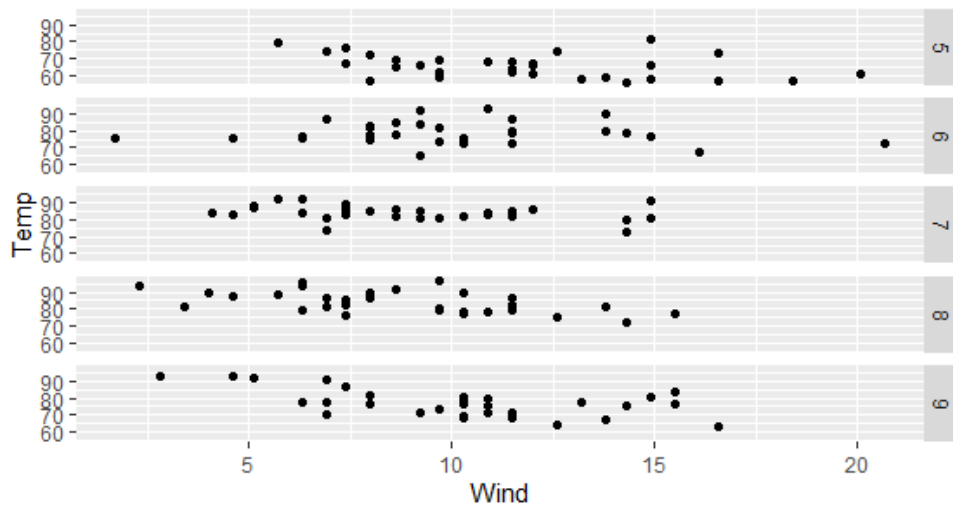
```

1 #按月份生成从左至右的多个散点图
2 qplot(Wind,Temp,data=airquality,facets = .~Month)

```

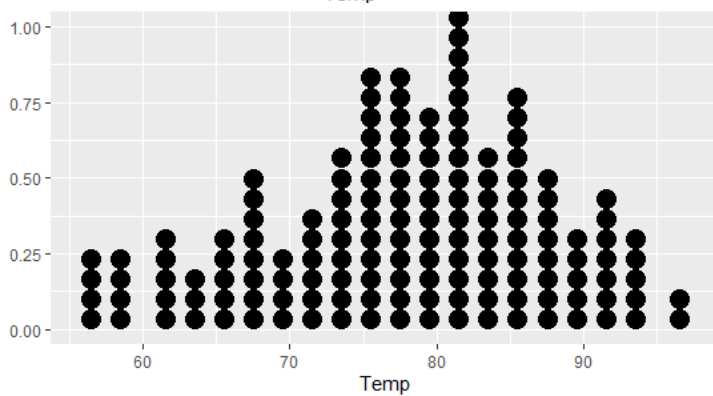
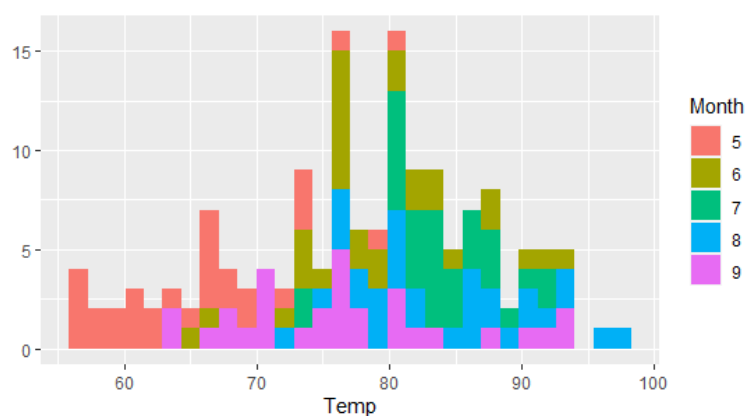
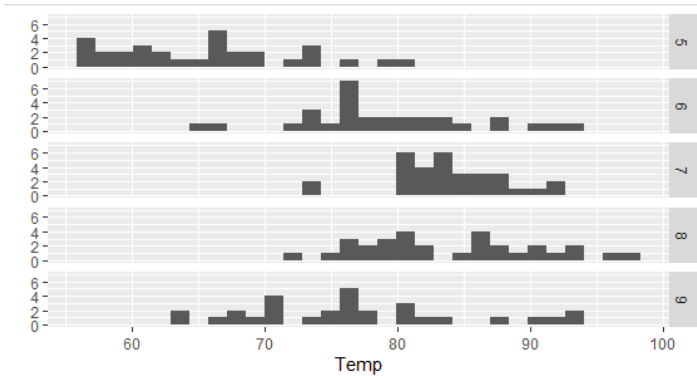
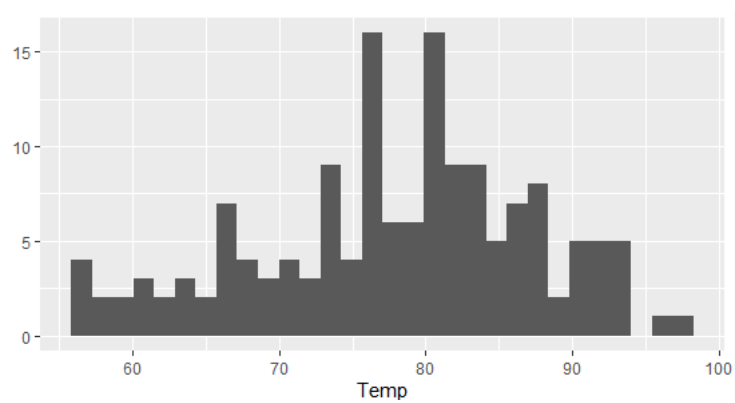


```
1 #按月份生成从上至下的多个散点图
2 qplot(Wind,Temp,data=airquality,facets = Month~.)
```

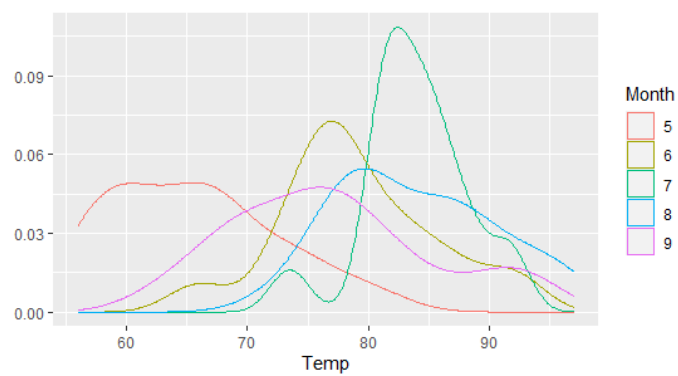
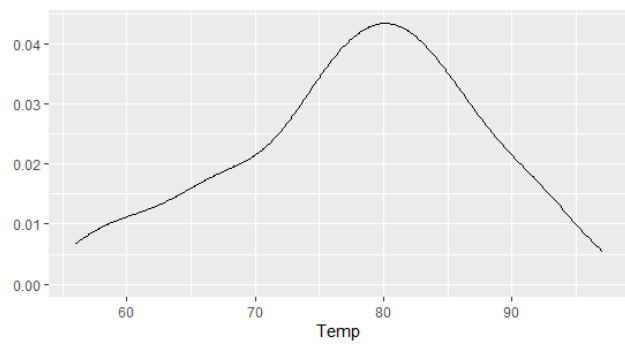


```
1 #可以看出，在定义参数facets时，若为~.变量名，则按变量不同从左到右生成不同的图
2 #若为变量名~.，则从上至下生成不同的图
```

```
1 #若在qplot函数中只输入了一个数据参数，则默认生成柱状图，如：
2 qplot(Temp,data=airquality)
3 qplot(Temp,data=airquality,facets = Month~.)
4
5 #在qplot函数中定义参数fill,可以得到一个累加的柱状图，如：
6 qplot(Temp,data=airquality,fill=Month)
7
8 #生成温度的密度点图
9 qplot(Temp,data=airquality,geom="dotplot")
10
```



```
1 #生成温度的密度曲线
2 qplot(Temp,data=airquality,geom = "density")
3 #按月份生成不同颜色的温度的密度曲线
4 qplot(Temp,data=airquality,geom = "density",color=Month)
```



Lattice绘图系统的绘图函数

lattice包

xyplot/bwplot/histogram/stripplot/dotplot/splom/leveplot/contourplot

- xyplot函数的格式: `xyplot(y~x|f*g,data)`,其中, `y`代表因变量, `x`代表自变量, `f*g`指分类变量, `data`指数据集, 第一个参数指一个公式, 其|的左半部分是必须存在的, 右半部分不是必须存在的, 若只存在左半部分, 说明我们不考虑交互作用, 只考虑`y`和`x`这两个变量之间的关系; 若左半部分和右半部分都有, 则考虑交互作用, 这样就提供给我们看了`x`和`y`这两者的关系在分类变量的不同水平下进行的变化。
- panel函数, 用于控制每个面板的绘图

grid包

- 实现了独立于base的绘图系统
- lattice包是基于grid创建的, 很少直接从grid包调用函数

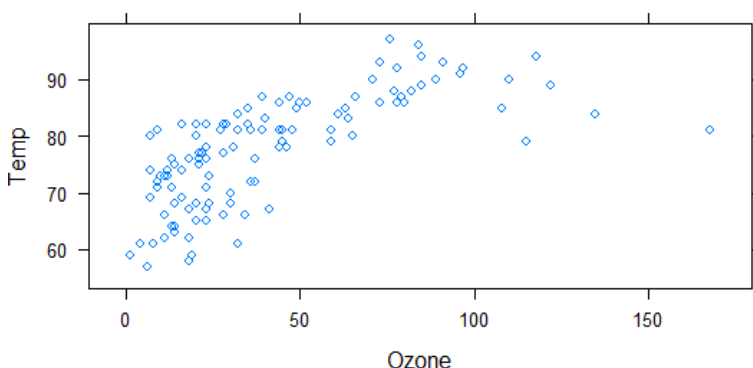
Lattice与Base的重要区别

Base绘图函数直接在图形设备上绘图

Lattice绘图函数返回trellis类对象

- 打印函数真正执行了在设备上绘图
- 命令执行时, trellis类对象会被自动打印, 所以看起来就像是lattice函数直接完成了绘图

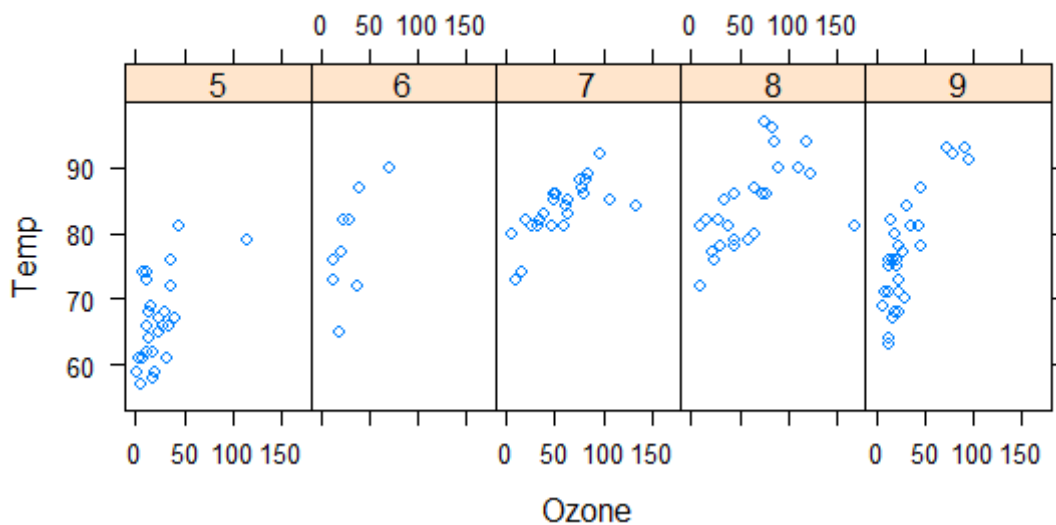
```
1 install.packages("lattice")
2 library(lattice)
3 #用xyplot()函数查看温度和臭氧含量的关系
4 xyplot(Temp~Ozone,data=airquality)
```



```

1 #将月份转化为分类变量
2 airquality$Month<-factor(airquality$Month)
3 airquality$Month
4 #查看不同月份下温度和臭氧含量的关系
5 xyplot(Temp~Ozone|Month,data=airquality,layout=c(5,1))
6 #结果显示：将月份转化为分类变量时有5个水平，分别为5、6、7、8、9月份，所以，我们在查看不同月份
7 #参数设置成了一行5列，其中，第一列指5月份温度和臭氧含量之间的散点图，第2列指6月份温度和臭氧含
8 #从总体上来看，温度和臭氧含量之间的关系随着月份的不同时有所变化的，这也意味着温度和臭氧含量之
9 #这也展示了lattice系统适合呈现交互作用的优势

```



Panel函数应用

```

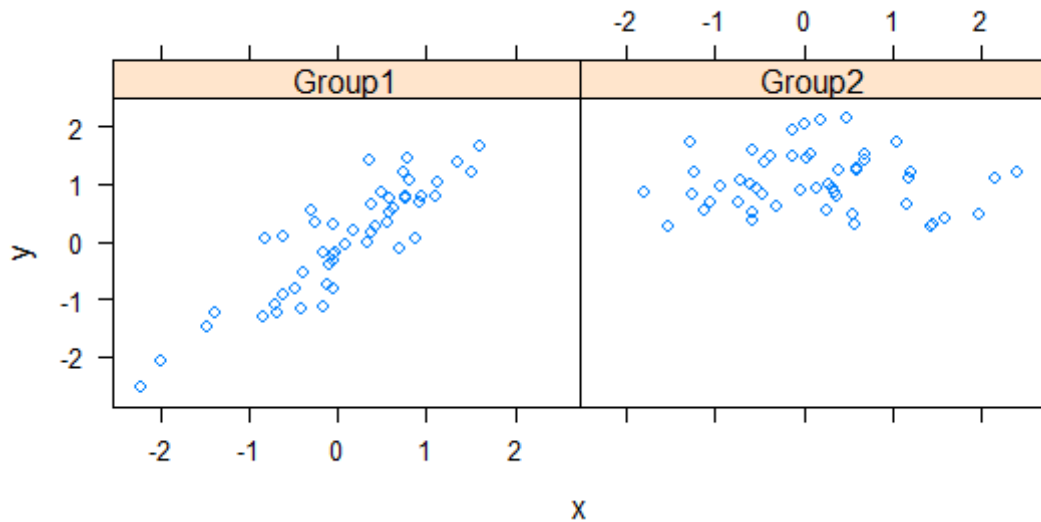
1 #1.用set.seed()函数设置种子数，其里面的数字可以输任何数，这里我们输入1
2 set.seed()
3 #设置种子点的意义是让我们每次产生的随机数是一样的
4
5 #2.在标准正态分布中产生100个随机数
6 x<-rnorm(100)
7 #3.重新创建一个变量f,该变量中只包含0和1这两个数，且每个数出现的次数均为50次
8 f<-rep(0:1,each=50)
9 #4.用x和f进行计算，将结果赋值给y
10 y<-x+f-f*x+rnorm(100,sd=0.5)
11 #目的是让x和y之间的关系与f有交互作用，我们为了让画出的图的点不在一条直线上，则加入一个随机数
12 #中产生100个随机数，设置正态分布的均值为0，标准差为0.5，否则正态分布的默认标准差为1
13
14
15 #5.把f变量转化为分类变量

```

```

16 f<-factor(f,labels = c("Group1","Group2"))#由于我们不知道变量f中的0和1代表什么，所以，我们
17
18 #6. 查看在f的不同水平下x和y之间的关系
19 xyplot(y~x|f,layout=c(2,1))

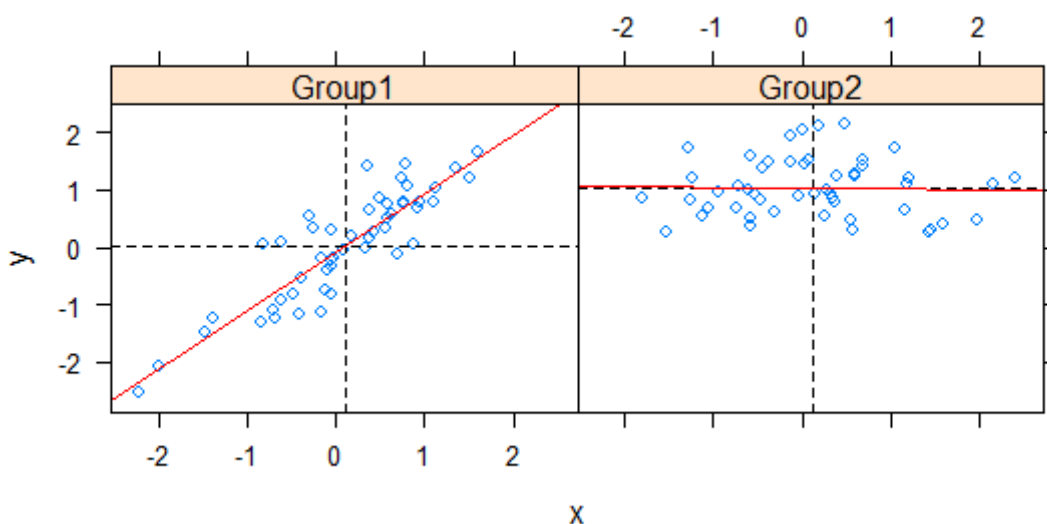
```



```

1 #7. 自己设置面板风格
2 xyplot(y~x|f,panel = function(x,y){panel.xyplot(x,y)
3   +panel.abline(v=mean(x),h=mean(y),lty=2)
4   +panel.lmline(x,y,col="red")})
5
6 #上述代码中，第二个参数指自己定义的函数，函数体用花括号括起来，panel.abline()函数中的第一个参
7 #panel.abline()函数中的第二个参数指在y的均值处画一条水平的直线，第三个参数的类型指线的类型是
8 #panel.lmline()函数的目的是在每个面板中添加一条水平的虚线和一条垂直的虚线，它们分别对应x的均
9 #这条红色的线是对数据进行拟合得到的回归线，在Group2水平下的结果和Group1类似

```



1.任意步长的等差数列构成的向量，用函数seq()

2.使用重复函数 rep()

```
1 rep(x,time=,each=length.out=)
2 #x:数量、向量或是数据对象，是要复制的对象
3 #times:表示x被复制的次数
4 #each:表示每个元素重复的次数
5 #length.out:表示截取前多少个元素
```

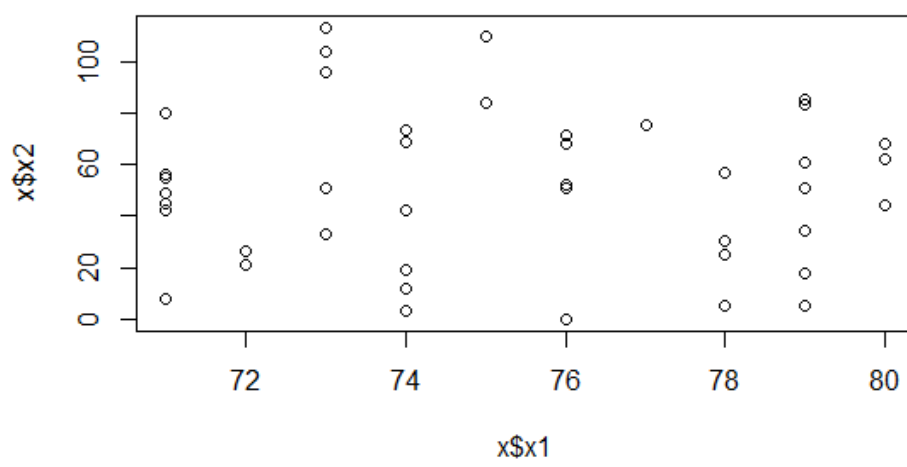
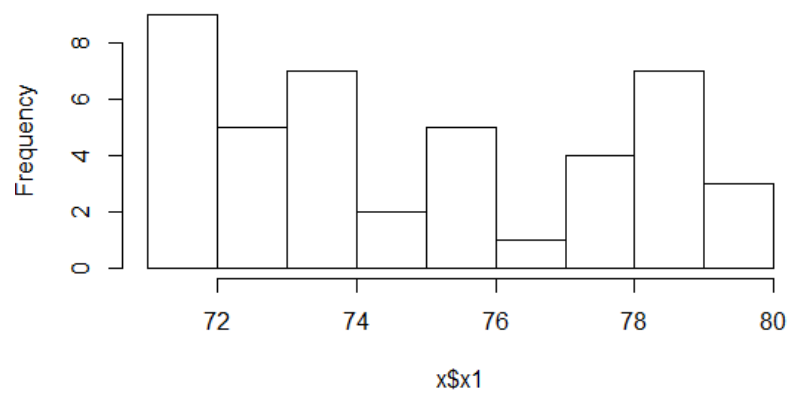
```
1 > #求余
2 > 9%4
3 [1] 1
4 > #求商
5 > 9/%4
6 [1] 2
7 > #乘方
8 > 2^3
9 [1] 8
10 > #开方
11 > sqrt(9)
12 [1] 3
```

例1：用一个班43位同学高数、线代、概率论与数理统计的成绩为例，在R语言中演示基本的统计分析

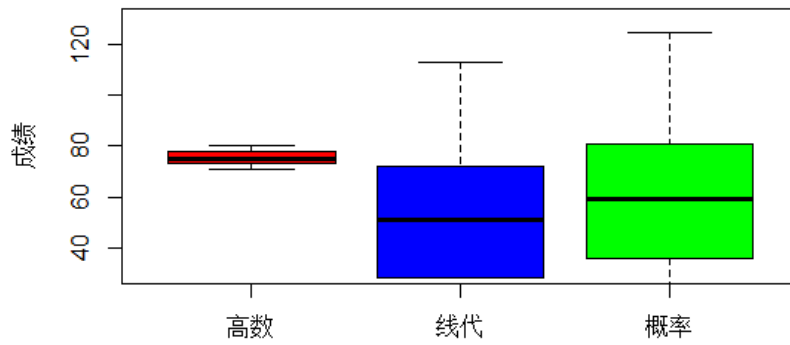
```
1 #1.生成43个学号
2 num<-seq(1010913418,1010913460)
3 num
4
5 #2.依次生成高数成绩x1,线代成绩x2,概率论与数理统计成绩x3为:
6 x1<-round(runif(43,min=70,max=80))
7 x2<-round(rnorm(43,mean=60,sd=30))
8 x3<-round(rnorm(43,mean=55,sd=29))
9 #round(x)表示对x取整数,runif(n,min=,max=)表示生成n个最小值为，最大值为 的服从均匀分布的数
10 #x2=round(rnorm(43,mean=60,sd=30))#表示生成43个均值为60，方差为30的服从正态分布的数
```

```
11
12 #3.将学号连同三科成绩合并成数据框x中
13 x<-data.frame(num,x1,x2,x3)
14 x
15
16 #4.计算各科的均值
17 colMeans(x[,2:4])
18 apply(x,2,mean)#matrix 1 indicates rows, 2 indicates columns, c(1, 2) indicates rows
19 apply(x[,2:4],2,mean)
20
21 #5.计算各科的最大、最小值
22 apply(x[,2:4],2,max)
23 apply(x[,2:4],2,min)
24
25 #6.求出每人的最高分
26 apply(x[,2:4],1,sum)
27
28 #7.求出最高分的同学程序为:
29 which.max(apply(x[,2:4],1,sum)) #which.max表示找到最高分同学的位置
30 x$num[which.max(apply(x[,2:4],1,sum))] #x$num表示找到最高分同学所对应的学号
31
32 #8.直方图
33 hist(x$x1)
34 #散点图
35 plot(x$x1,x$x2)
36 boxplot(x[2:4],main="学生成绩箱线图",ylab="成绩",names=c("高数","线代","概率"),col=c("red","blue","green"))
37 plot(x$x1,x$x2,main="高数成绩与线代成绩散点图",xlab="高数成绩",ylab="线性代数成绩",xlim=c(0,100),ylim=c(0,100))
```

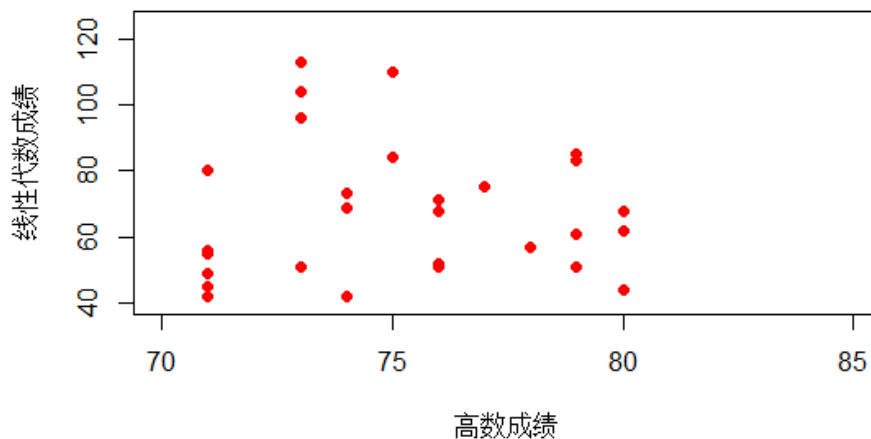
Histogram of x\$x1



学生成绩箱线图



高数成绩与线代成绩散点图



```

1 #常见的绘图参数解释:
2 #font=字体;lty=线类型;lwd=线宽度;pch=点的类型;xlab=横坐标;ylab=纵坐标;xlim=横坐标范围
3
4 #coplot(x~y|z) 在z的每个值或每个区间上做x与y的散点图
5 #pairs(x) 当x为矩阵,做x各列之间的散点图
6 #hist(x,freq)直方图,参数freq默认为TRUE,根据频数作图;若为FALSE,则根据构成比(总和为1)作
7 #barplot(table(x,y),beside=FALSE) 对定性变量x,y做条图,默认堆积条图,使用beside=T则为并列
8 #qqnorm(x) QQ图,(正态分位数-分位数图)
9 #pie(table(x))#对定性变量x做饼图
10 #lines(x),lines(x,y) 添加折线原图上修改
11 #abline(lm(y~x)) 添加y对x的回归直线
12 #abline(a,b) a为截距,b为斜率
13 #abline(v=) 添加垂直线
14 #abline(h=) 添加水平线
    
```

```

1 #使用sort(colors())可以查看所有已知的颜色
    
```


假设检验 (hypothesis test) 亦称显著性检验 (significant test),是统计推断的另一个重要内容,其目的是比较总体参数之间有无差别。假设检验的实质是判断观察到的“差别”是由抽样误差引起还是总体上的不同,目的是评价两种不同处理引起效应的不同的证据有多强,这种证据的强度用概率P来度量和表示。除t分布外,针对不同的资料还有其他各种检验统计量及分布,如F分布、 χ^2 分布等,应用这些分布对不同类型的数据进行假设检验的步骤相同,其差别仅仅是需要计算的检验统计量不同。

其基本思想是小概率反证法思想

1.W检验(Shapiro-Wilk)

检验数据是否符合正态分布

R函数: shapiro.test()

结果含义: 当p值小于某个显著性水平 α (比如0.05) 时, 则认为样本不是来自正态分布的总体, 否则承认样本来自正态分布的总体

2.K检验(经验分布的Kolmogorov-Smirnov检验)

R函数: ks.test(),如果P值很小, 说明拒绝原假设表明数据不符合F(n,m)分布。

3.相关性检验:

R函数: cor.test()

```
1 cor.test(x,y,alternative=c("two.sided","less","greater"),method=c("pearson","kendall")
```

结果含义: 如果p值很小, 则拒绝原假设, 认为x,y是相关的。否则认为是不相关的。

4.T检验

用于正态总体均值假设检验, 单样本、双样本都行

```
1 t.test()  
2 t.test(x,y=NULL,alternative=c("two.sided","less","greater"),mu=0,paired=FALSE,var.equ
```

结果含义: P值小于显著性水平时拒绝原假设, 否则, 接受原假设。具体的假设要看所选择的是双边假设还是单边假设 (又分小于和大于)

5.Fisher精确的独立检验

原假设：X, Y相关

```
1 fisher.test(x,y=NULL,workspace=200000,hybrid=FALSE,control=list(),or=1,alternative="t
```

6.Pearson χ^2 拟合优度检验

R语言中调用chisq.test(X)

原假设H0：X符合F分布

p-值小于某个显著性水平，则表示拒绝原假设，否则接受原假设

正态分布参数检验

例1.某种原件的寿命X(以小时计)服从正态分布N (μ , σ)。其测得16只原件的寿命如下：

159 280 101 212 224 379 179 264

222 362 168 250 149 260 485 170

问是否有理由认为元件的平均寿命大于255小时？

解：按题意，需检验H0： $\mu \leq 255$, H1 > 255

此问题属于单边检验问题，可以使用R语言t.test(x,y=NULL,#只提供x为单个正态总体均值检验，否则为两个总体均值检验

alternative=c("two.side","less","greater"),#双边检验 单边检验

mu=0,#原假设：mu=0,均值为某个具体数字

paired=FALSE,

var.equal=FALSE,# 方差齐性选项

conf.level=.95)#置信水平95%

```
1 > X<-c(159,280,101,212,224,379,179,264,222,362,168,250,149,260)
2 > t.test(X,alternative = "greater",mu=225)
3
4     One Sample t-test
5
6 data:  X
7 t = 0.20141, df = 13, p-value = 0.4217
```

```

8 alternative hypothesis: true mean is greater than 225
9 95 percent confidence interval:
10 192.1588      Inf
11 sample estimates:
12 mean of x
13 229.2143

```

可见P值为0.257>0.05,不能拒绝原假设, 接受H0, 即平均寿命不大于225小时

例2.在平炉上进行的一项实验以确定改变操作方法的建议是否会增加钢的得率, 实验时在同一个平炉上进行的, 每炼一炉钢时除操作方法外, 其他条件都尽可能做到相同, 先用标准方法炼一炉, 然后用新办法炼一炉, 以后交替进行, 各炼了10炉, 其得率分别为标准方法: 78.1 72.4 76.2 74.3 77.4 78.4 76.0 75.5 76.7 77.3

新方法: 79.1 81 77.3 79.1 80 79.1 79.1 77.3 80.2 82.1

设这两个样本相互独立, 且分别来自正态总体 $N(\mu_1, \sigma_1)$ 和 $N(\mu_2, \sigma_2)$, 其中 μ_1 , μ_2 和 σ_2 未知。问新的操作能否提高得率? (取 $\alpha=0.05$)

解: 因为数据是成对出现的, 所以采用成对数据t检验比上述的双样本均值检验更准确。所谓成对t检验就是 $Z_i = X_i - Y_i$, 再对Z进行单样本均值检验

```

1 > X<-c(78.1,72.4,76.2,74.3,77.4,78.4,76.0,75.5,76.7,77.3)
2 > Y<-c(79.1,81.0,77.3,79.1,80.0,79.1,77.3,80.2,82.1)
3 > t.test(X-Y,alternative = "less")
4
5 #结果如下:
6     One Sample t-test
7
8 data:  X - Y
9 t = -3.9114, df = 9, p-value = 0.001779
10 alternative hypothesis: true mean is less than 0
11 95 percent confidence interval:
12      -Inf -1.700279
13 sample estimates:
14 mean of x
15      -3.2

```

可见P值<0.05,接受备择假设, 即新操作能够提高得率。并且P值更小可见比双样本均值检验更准确。

例3.对例2进行方差检验，方差是否相同

```
1 > var.test(X,Y,ratio = 1,alternative = c("two.sided","less","greater"),conf.level = 0.05)
2
3 #结果如下:
4      F test to compare two variances
5
6 data:  X and Y
7 F = 1.3365, num df = 9, denom df = 8, p-value = 0.6933
8 alternative hypothesis: true ratio of variances is not equal to 1
9 95 percent confidence interval:
10  0.3067324 5.4822833
11 sample estimates:
12 ratio of variances
13      1.336505
```

可见P值为0.559>0.05,接受原假设，认为两者方差相同

例4.皮尔森拟合优度检验

某消费者协会为了确定市场上消费者对5种品牌啤酒的喜好情况，随机抽取了1000名啤酒爱好者作为样品进行试验：每个人得到5种品牌的啤酒各一瓶，但未表明牌子。这5种啤酒分别按着A、B、C、D、E字母的5张纸片随即的顺序送给每一个人。下表是根据样本资料整理的各种品牌啤酒爱好者的频数分布。试根据这些数据判断消费者对这5种品牌啤酒的爱好有无明显差异？

最喜欢的牌子	A	B	C	D	E
人数X	210	312	170	85	223

解：如果消费者对5种品牌的啤酒无显著差异，那么，就可以认为喜好这5种品牌啤酒的人呈均匀分布，即5种品牌啤酒爱好者人数各占20%。据此假设

H0：喜好5种啤酒的人数分布均匀

可以使用Pearson χ^2 拟合优度检验，R语言中调用chisq.test(X)

```
1 > X<-c(210,312,170,85,223)
2 > chisq.test(X)
3
4      Chi-squared test for given probabilities
```

```
5
6 data: X
7 X-squared = 136.49, df = 4, p-value < 2.2e-16
```

例5.正态W检验

已知15名学生体重如下，问是否服从正态分布 75 64 47.4 66.9 62.2 62.2 58.7 63.5 66.6 64 57 69 50 72

```
1 > w<-c(75,64,47.4,66.9,62.2,62.2,58.7,63.5,66.6,64,57,69,50,72)
2 > shapiro.test(w)
3
4      Shapiro-Wilk normality test
5
6 data:  w
7 W = 0.95448, p-value = 0.6321
```

因为P值>0.05,接受原假设，所以认为来自正态分布总体

当我们在R中读入了一个txt\csv等格式的数据时，即自成了一个数据框

数据框就是一个表格，它的每一列可以存放不同类型的数据，因此，它最接近于实际数据的形态

一、创建数据框

```
1 > birthday<-c(1997,1998,1971,1999,2000,2001,2003)
2 > gender<-c("男","男","女","男","女","女","男")
3 > people<-data.frame(name,birthday,gender);people
4   name birthday gender
5 1  张三     1997     男
6 2  李四     1998     男
7 3  王大妈   1971     女
8 4   小二     1999     男
9 5   小五     2000     女
10 6   小六     2001     女
11 7   小七     2003     男
```

二、筛选、引用

```
1 > people[3,1]
2 [1] 王大妈
3 Levels: 李四 王大妈 小二 小六 小七 小五 张三
4 > people[3,1]#查看第3行第1列的人物信息
5 [1] 王大妈
6 Levels: 李四 王大妈 小二 小六 小七 小五 张三
7 > people[3,] #查看第3行的人物信息
8   name birthday gender
9 3 王大妈     1971     女
10 > people$birthday
11 [1] 1997 1998 1971 1999 2000 2001 2003
12 > people[,2]#查看第2列人物的出生年份
13 [1] 1997 1998 1971 1999 2000 2001 2003
14 > new_shuju<-people[people$birthday==2000,];new_shuju
15   name birthday gender
16 5 小五     2000     女
17 > #选择出生年份为2000年的数据
18
19 > new_shuju<-people[people$gender=="男"&people$birthday>1998,]
```



```

20 > new_shuju
21   name birthday gender
22 4 小二      1999     男
23 7 小七      2003     男

```

三、增列、合并

```

1 #1.增列
2 > nation<-c("汉族","汉族","回族","藏族","汉族","汉族","苗族")
3 > #给数据框添加一行数据：民族
4 > people$nat<-nation
5 > people$nation<-nation
6 > #给数据框的增列命名了2次，因此要删除其中的一列
7 > people[, -5]
8   name birthday gender  nat
9 1 张三      1997     男 汉族
10 2 李四      1998     男 汉族
11 3 王大妈    1971     女 回族
12 4 小二      1999     男 藏族
13 5 小五      2000     女 汉族
14 6 小六      2001     女 汉族
15 7 小七      2003     男 苗族

```

```

1 #2.合并
2 > x1<-people[1:3,];x1
3   name birthday gender  nat nation
4 1 张三      1997     男 汉族  汉族
5 2 李四      1998     男 汉族  汉族
6 3 王大妈    1971     女 回族  回族
7 > x2<-people[4:7,];x2
8   name birthday gender  nat nation
9 4 小二      1999     男 藏族  藏族
10 5 小五      2000     女 汉族  汉族
11 6 小六      2001     女 汉族  汉族
12 7 小七      2003     男 苗族  苗族
13 > people.nat<-merge(x1,x2,by="nat");people.nat
14   nat name.x birthday.x gender.x nation.x name.y birthday.y gender.y nation.y
15 1 汉族 张三      1997     男    汉族 小五      2000     女    汉族

```

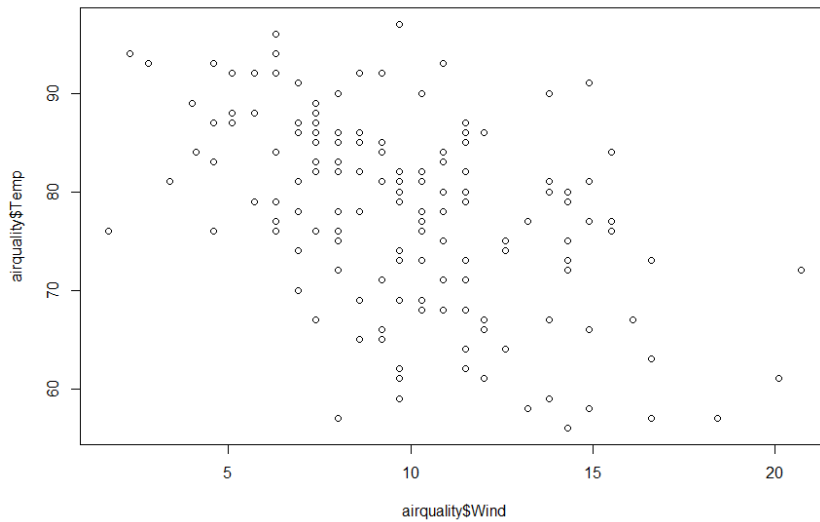
16	2	汉族	张三	1997	男	汉族	小六	2001	女	汉族
17	3	汉族	李四	1998	男	汉族	小五	2000	女	汉族
18	4	汉族	李四	1998	男	汉族	小六	2001	女	汉族

```

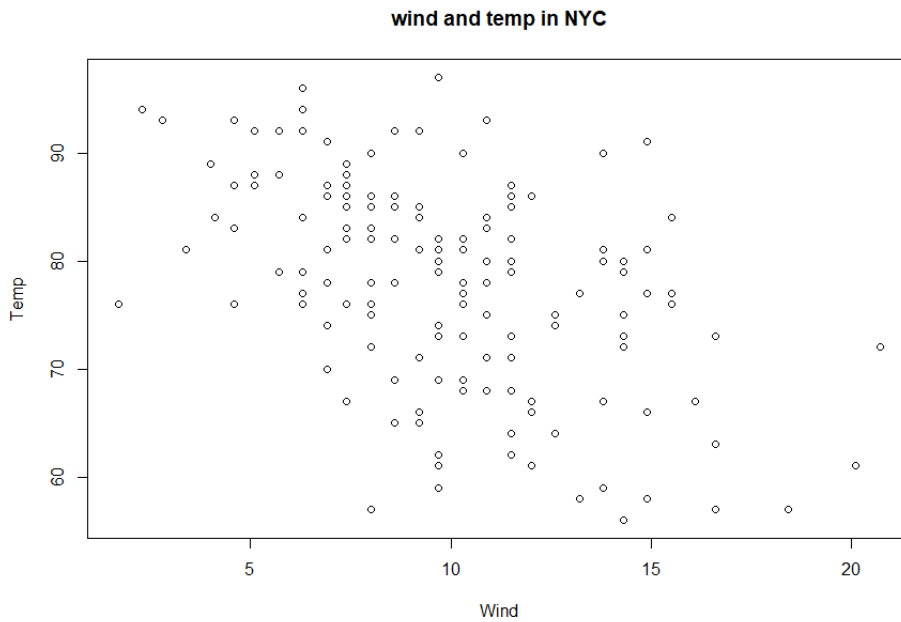
1 #显示匹配不到的样本，只需要在merge()函数中多传入一个参数all.x=T即可
2 > #all.x=T,将前一个数据框people中nat列所有的值做合并，匹配不到赋值NA
3 > x3<-x2[1:3,];x3
4   name birthday gender  nat nation
5 4 小二      1999     男 藏族   藏族
6 5 小五      2000     女 汉族   汉族
7 6 小六      2001     女 汉族   汉族
8
9 > people.nat<-merge(x1,x3,by="nat",all.x = T);people.nat
10    nat name.x birthday.x gender.x nation.x name.y birthday.y gender.y nation.y
11 1 汉族   张三      1997     男   汉族   小五      2000     女   汉族
12 2 汉族   张三      1997     男   汉族   小六      2001     女   汉族
13 3 汉族   李四      1998     男   汉族   小五      2000     女   汉族
14 4 汉族   李四      1998     男   汉族   小六      2001     女   汉族
15 5 回族 王大妈      1971     女   回族   <NA>      NA      <NA>   <NA>

```

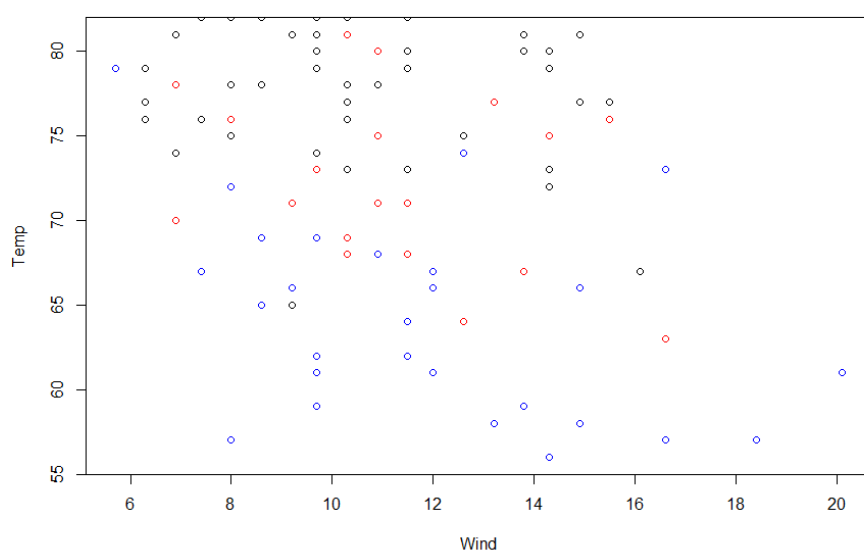
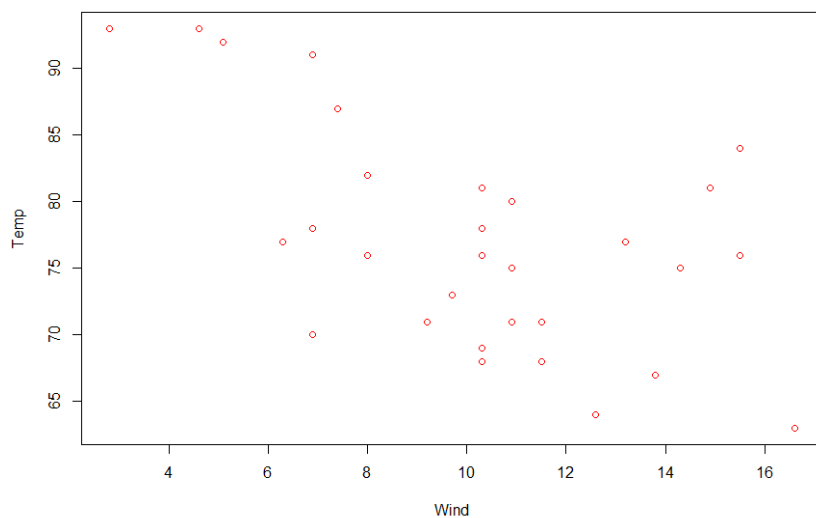
```
1 #用plot()函数画空气质量这个数据集里风速和温度的散点图
2 plot(airquality$Wind,airquality$Temp)
```



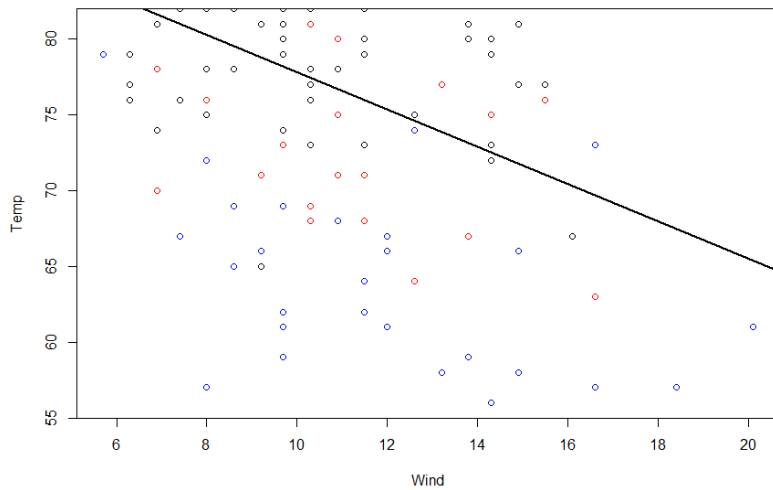
```
1 #用with()函数画空气质量这个数据集里风速和温度的散点图
2 with(airquality,plot(Wind,Temp))
3 #with()函数的第一个参数为数据集，第二个参数指绘图函数
4
5 #给散点图添加标题
6 #法1
7 with(airquality,plot(Wind,Temp,main="wind and temp in NYC"))
8 #法2
9 with(airquality,plot(Wind,Temp))
10 title(main="wind and temp in NYC")
```



```
1 #按月份来画点，不同月份对应的数值显示不同的颜色
2 #一、画空气质量数据集中风速和温度9月份的散点图
3 #1.筛选空气质量数据集的子集，子集是9月份的数据
4 x<-subset(airquality,Month==9)
5
6 #2.画图
7 with(x,points(Wind,Temp,col="red"))
8 #画空气质量数据集中风速和温度5月份的散点图
9 with(subset(airquality,Month==5),points(Wind,Temp,col="blue"))
10 #画空气质量数据集中风速和温度6、7、8月份的散点图
11 with(subset(airquality,Month%iin% c(6,7,8)),points(Wind,Temp,col="black"))
```



- 1 #在图形中添加回归线
- 2 #1.先用lm()函数拟合一个线性模型
- 3 `fit<-lm(Temp~Wind,airquality)`
- 4 #2.给图形添加回归线，并设置线宽
- 5 `abline(fit,lwd=2)`
- 6 #上述lm()函数拟合的线性模型中，Temp(风速)指因变量，Wind(温度)指自变量，它们都来自数据集airq



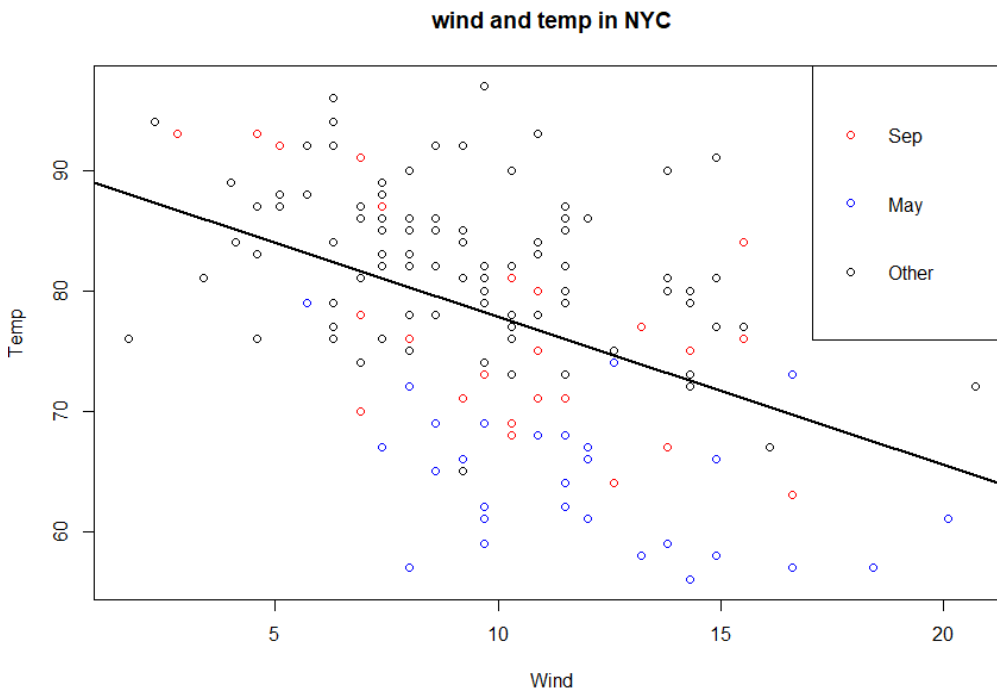
1 #五、用legend()函数给图形添加图例

2 legend("topright",pch=1,col=c("red","blue","black"),legend=c("Sep","May","Other"))

3 #legend()函数的第一个参数指图例说明的位置在右上方，第二个参数指图例保持跟散点图一样的蓝色、黑

4 #第三个参数指颜色，对应于我们画图的顺序，第四个参数指颜色赋予的含义，Others指6/7/8月份

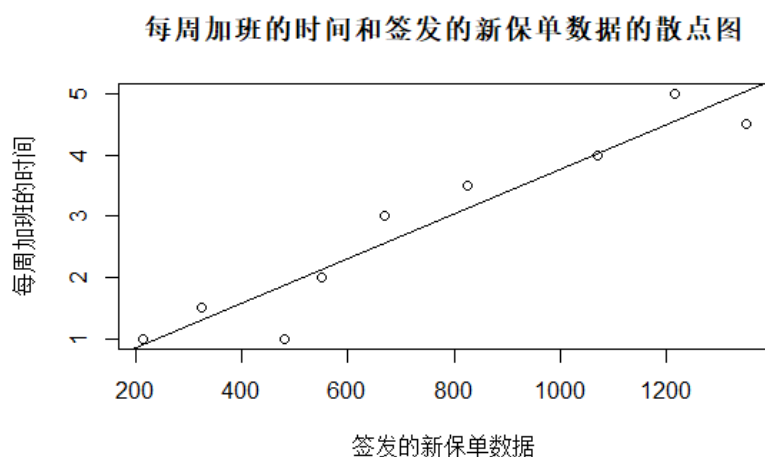
5 legend("topleft",pch=1,col=c("red","blue","black"),legend=c("Sep","May","Other"))



[ggplot绘图系统之qplot](#)

回归分析 (Regression Analysis) 是用来确定2个或2个以上变量间关系的一种统计分析方法。如果回归分析中, 只包括了一个自变量X和一个因变量Y时, 且它们的关系是线性的, 那么这种回归分析称为一元线性回归分析。

```
1 x<-c(825,215,1070,550,480,1350,325,670,1215)
2 y<-c(3.5,1,4,2,1,4.5,1.5,3,5)
3 library(ggplot2)
4 plot(x,y,main="每周加班的时间和签发的新保单数据的散点图",xlab="签发的新保单数据",ylab="每周
5 abline(lm(y~x))
6
7 #每周加班的时间和签发的新保单数据呈线性关系, 因而可以考虑一元线性回归模型
```



```
1 #2. 求回归方程, 并对相应的方程做检验。程序中lm()函数表示作线性模型, summary()函数提取模型的计
2 a<-lm(y~x)
3 > summary(a)
4
5 Call:
6 lm(formula = y ~ x)
7
8 Residuals:
9      Min       1Q   Median       3Q      Max
10 -0.87004 -0.12503  0.09527  0.37323  0.45258
11
12 Coefficients:
13      Estimate Std. Error t value Pr(>|t|)
```



```

14 (Intercept) 0.1215500 0.3588377 0.339 0.745
15 x          0.0036427 0.0004303 8.465 6.34e-05 ***
16 ---
17 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
18
19 Residual standard error: 0.485 on 7 degrees of freedom
20 Multiple R-squared:  0.911, Adjusted R-squared:  0.8983
21 F-statistic: 71.65 on 1 and 7 DF, p-value: 6.344e-05

```

残差的标准差 (Residual standard error) 为0.485，其自由度为 7 (计算方式: $n-2$);

相关系数的平方 (Multiple R-squared) 即判别系数为 0.911,

Adjusted R-squared (调整后的R方): 0.8983，说明整个模型的拟合度较高，两者存在较强的线性关系

F统计量(F-statistic) 的自由度(1,7),回归方程为: **Intercept的意思是截距**

$$\hat{y} = 0.12155 + 0.0036427x$$

```

1 #3.参数的区间估计
2 > confint(a)
3           2.5 %      97.5 %
4 (Intercept) -0.726966290 0.970066254
5 x           0.002625117 0.004660271

```

常数项的置信区间为[-0.726966290,0.970066254],系数的置信区间为 [0.002625117,0.004660271]

```

1 #4.预测
2 > new<-data.frame(x=600)
3 > new
4     x
5 1 600
6
7 > new2<-predict(a,new,interval = "confidence",level = 0.95)
8 > new2
9      fit      lwr      upr
10 1 2.307166 1.897614 2.716719

```

平均值的置信区间为[1.897614,2.716719]

```
1 #5.残差分析
2 > e<-resid(a,digits=5) #计算回归a的残差，并保留5位小数
3 > e
4           1           2           3           4           5           6
5 0.37322742 0.09527080 -0.01923262 -0.12503171 -0.87004313 -0.53918695
6           7           8           9
7 0.19457445 0.43784500 0.45257674
```

```
1 library(deplyr)
2 library(psych) # 本例中用到describe函数(描述统计)
3 library(lm.beta) # 用于在回归分析中输出标准化回归系数 $\beta$ 
4 library(ggplot2) # 绘图必备
5 library(gridExtra) # 多图组合
6 library(PerformanceAnalytics) # 本例中用到chart.Correlation函数
7
8 #描述统计
9 data=get(data("anscombe")) # 统计学经典四组数据，具有相同的统计量(M, SD, r)，但本质不同
10 describe(data)
11 chart.Correlation(data) # 绘制变量分布及相关关系图
12
13 # 线性回归：以下4种写法在作用上等价(建模、赋值、报告结果) -----
14 # y1 = b0 + b1*x1
15 model.1=lm(y1 ~ x1, data)
16 summary(lm.beta(model.1))
17 # y2 = b0 + b1*x2
18 model.2=lm(y2 ~ x2, data)
19 model.2 %>% lm.beta() %>% summary()
20 # y3 = b0 + b1*x3
21 summary(lm.beta( (model.3=lm(y3 ~ x3, data)) ))
22 # y4 = b0 + b1*x4
23 (model.4=lm(y4 ~ x4, data)) %>% lm.beta() %>% summary()
24
25 # 散点图 -----
26 p1=ggplot(data=data, aes(x=x1, y=y1)) + geom_point() + geom_smooth(method="lm") + lab
27 p2=ggplot(data=data, aes(x=x2, y=y2)) + geom_point() + geom_smooth(method="lm") + lab
28 p3=ggplot(data=data, aes(x=x3, y=y3)) + geom_point() + geom_smooth(method="lm") + lab
29 p4=ggplot(data=data, aes(x=x4, y=y4)) + geom_point() + geom_smooth(method="lm") + lab
30 grid=grid.arrange(grobs=list(p1, p2, p3, p4), ncol=2, nrow=2)
31 ggsave("Plot.pdf", grid, width=8, height=6) # PDF保存的是矢量图，任意放缩不会模糊，推荐使
32 ggsave("Plot.png", grid, width=8, height=6, dpi=300) # 科研绘图一般要求分辨率至少达到 300
```

factoextra是一个R软件包，可以轻松提取和可视化探索性多变量数据分析的输出，其中包括：主成分分析（PCA），用于通过在不丢失重要信息的情况下减少数据的维度来总结连续（即定量）多变量数据中包含的信息。对应分析（CA），它是适用于分析由两个定性变量（或分类数据）形成的大型列联表的主成分分析的扩展。多重对应分析（MCA），它是将CA改编为包含两个以上分类变量的数据表格。多因素分析（MFA）专用于数据集，其中变量按组（定性和/或定量变量）组织。分层多因素分析（HMFA）：在数据组织为分层结构的情况下，MFA的扩展。