

# R基础统计及线性模型

陈华珊 (中国社科院社会发展战略研究院)

# 基础统计

# ◆ 描述统计

## ❖ 向量

---

<code>mean(x)</code>	Mean of x
<code>median(x)</code>	Median of x
<code>var(x)</code>	Variance of x
<code>sd(x)</code>	Standard deviation of x
<code>cov(x,y)</code>	Covariance of x and y
<code>cor(x,y)</code>	Correlation of x and y
<code>min(x)</code>	Minimum of x
<code>max(x)</code>	Maximum of x
<code>range(x)</code>	Range of x
<code>quantile(x)</code>	Quantiles of x for the given probabilities

---

```
1 x <- rnorm(10)
2 y <- runif(10)
3 mean(x)
```

```
4    ## [1] -0.2526202
5    cor(x,y)
6    ## [1] 0.2570403
```

## ❖ data.frame的变量 (列)

---

<code>summary()</code>	对于每列, 显示基本信息
<code>apply()</code>	对于每列, 执行一个函数
<code>tapply()</code>	拆分为多个子集, 执行一个函数, 返回数组
<code>by()</code>	类似于 <code>tapply()</code> `, 返回`by` 类
<code>ave()</code>	类似于 <code>`tapply()'</code> , 返回向量
<code>aggregate()</code>	类似于 <code>`tapply()'</code> , 返回 <code>dataframe</code>

---

## ❖ 示例

```
1  # Basic summary statistics for each column
2  summary(gss)
3
4  # Variance of first two columns
5  apply(gss[,1:2], 2, var)
6
7  # Mean family income by region
8  tapply(gss$a8a, gss$s41, FUN=mean) # Return a vector
9  by(gss$a8a, gss$s41, FUN=mean) # Return a by object,
10
11 # extract elements using []
12 ave(gss$a8a, gss$s41, FUN=mean) # Return a vector, same
13
14 # length as first argument
15 aggregate(gss$a8a, list(gss$s41) , FUN=mean) # Return a dataframe,
   subset
```

```
16
17   # variable needs to be a list
18   # Mean inc by region and area
19   with(gss, tapply(a8a, list(s41, s5a), FUN=mean))
20   with(gss, by(a8a, list(s41, s5a), FUN=mean))
21   with(gss, ave(a8a, list(s41, s5a), FUN=mean))
22   with(gss, aggregate(a8a, list(s41, s5a), FUN=mean))
```



## ❖ 使用第三方包

- `Hmisc::describe()` 返回观测数量、缺失值和唯一值的数目、平均值、分位数, 以及五个最大值和最小值;
- `pastecs::stat.desc()` 返回所有值、空值、缺失值的数量, 最小值、最大值、值域、标准差、均值, 均值95%置信区间等;
- `psych::describe()`
- `doBy::summaryBy()`

## ◆ 频次表/列联表

## ❖ 相关函数

---

<code>table(var1, var2, ..., varN)</code>	使用N个类别变量 (因子) 创建一个N维列联表
<code>xtabs(formula, data)</code>	根据一个公式和一个矩阵/数据框创建一个N维列联表
<code>prop.table(table, margins)</code>	百分比列联表, 行/列
<code>margin.table(table, margins)</code>	边际求和
<code>addmargins(table, margins)</code>	包含边际求和
<code>ftable(table)</code>	“平铺式” 列联表

---

## ❖ 示例

VCD::Arthritis数据集: 关于一项风湿性关节炎新疗法的双盲临床试验结果。

```
1 library(vcd)
2 mytable = with(Arthritis, table(Improved))
3 mytable
4 ## Improved
5 ##      None      Some Marked
6 ##      42      14      28
7
8 # 转化为百分比
9 prop.table(mytable)
10 ## Improved
11 ##      None      Some      Marked
12 ## 0.5000000 0.1666667 0.3333333
13
14 # 二维列联表
15 mytable = xtabs(~Treatment + Improved, data = Arthritis)
```

```

16 mytable
17 ##           Improved
18 ## Treatment None Some Marked
19 ##   Placebo   29   7   7
20 ##   Treated   13   7  21
21 margin.table(mytable, 1)
22 ## Treatment
23 ## Placebo Treated
24 ##      43      41
25 prop.table(mytable, 1)
26 ##           Improved
27 ## Treatment      None      Some      Marked
28 ##   Placebo 0.6744186 0.1627907 0.1627907
29 ##   Treated 0.3170732 0.1707317 0.5121951
30 addmargins(mytable)
31 ##           Improved
32 ## Treatment None Some Marked Sum
33 ##   Placebo   29   7   7  43

```

```

34  ##    Treated    13     7     21    41
35  ##    Sum       42    14     28    84
36  addmargins(prop.table(mytable))
37  ##              Improved
38  ## Treatment      None      Some      Marked      Sum
39  ##   Placebo 0.34523810 0.08333333 0.08333333 0.51190476
40  ##   Treated 0.15476190 0.08333333 0.25000000 0.48809524
41  ##    Sum    0.50000000 0.16666667 0.33333333 1.00000000
42  ###useNA = 'ifany'

```

## ◆ 均值比较

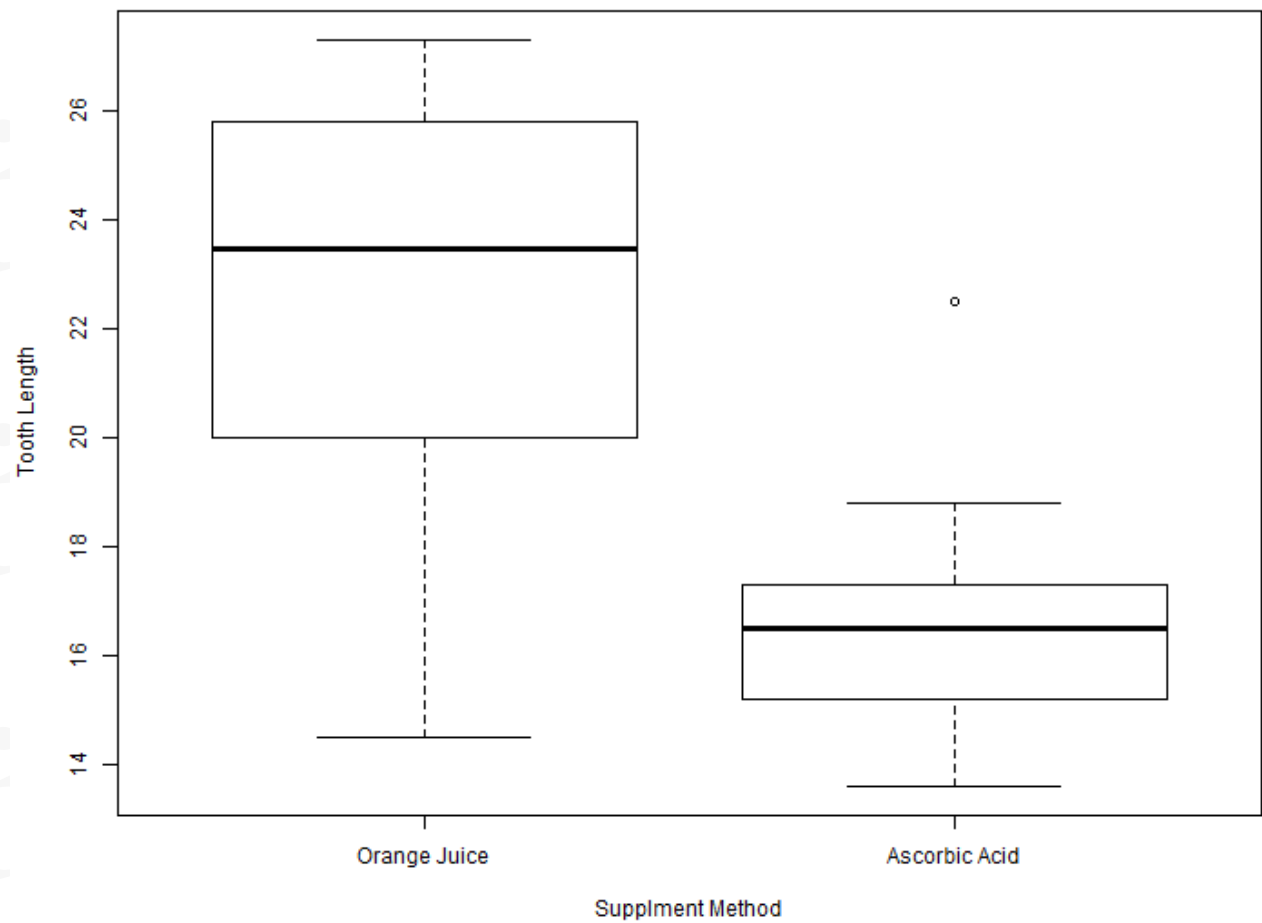
$t$ -test



```
t.test()
```

## ❖ 示例

```
1  # Only focus on the 1mg dose
2  Tooth.1mg <- subset(ToothGrowth, dose==1)
3  # From the boxplot orange juice leads to longer teeth
4  with(Tooth.1mg, boxplot(len~supp, names=c("Orange Juice", "Ascorbic
   Acid"),
5                                     xlab="Supplment Method", ylab="Tooth Length"))
```



```
1  # Two-sample t-test
2  tt <- t.test(len~supp, data=Tooth.1mg, alternative="two.sided",
3              var.equal=FALSE, conf.level=.95)
4  # Extract results
5  # See the documentation for a description of the values returned
6  names(tt)
7  ## [1] "statistic" "parameter" "p.value" "conf.int" "estimate"
8  ## [6] "null.value" "alternative" "method" "data.name"
9  tt$p.value
10 ## [1] 0.001038376
11 tt$conf.int
12 ## [1] 2.802148 9.057852
13 ## attr(,"conf.level")
14 ## [1] 0.95
```

## ◆ 方差比较

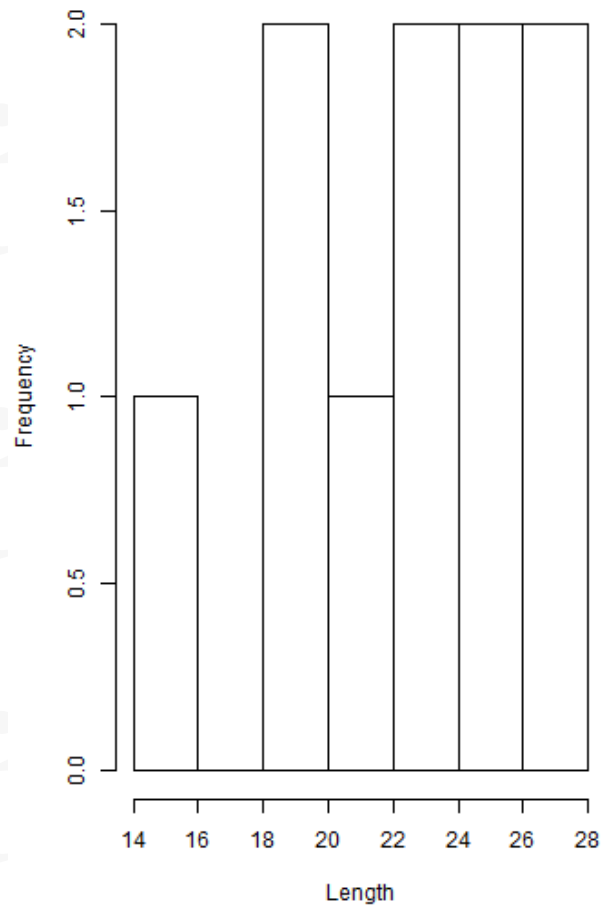
对于来自正态总体的两个样本，可用  $F$ -test 进行方差检验

```
1  `var.test()``
```

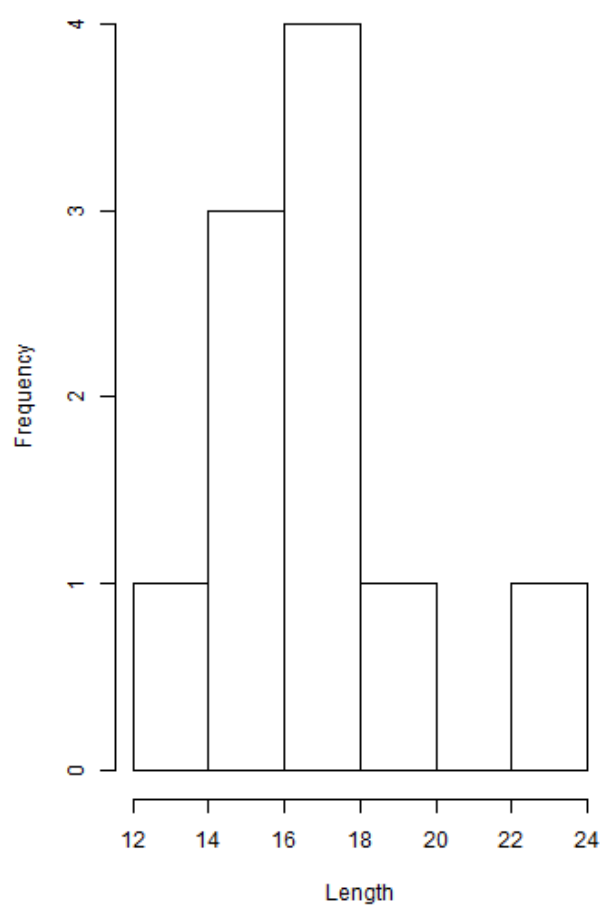
## ❖ 示例

```
1  # ToothGrowth dataset where dose=1mg
2  Tooth.1mg <- subset(ToothGrowth, dose==1)
3  # Sample Variances
4  with(Tooth.1mg, tapply(len, supp, var))
5  ##          OJ          VC
6  ## 15.295556  6.326778
7  # From histograms the normality assumption of the F-test is not
   reasonable
8  par(mfrow=c(1,2))
9  with(Tooth.1mg, hist(len[supp=="OJ"],
10                        main="Histogram of length for Orange Juice",
11                        xlab="Length"))
11 with(Tooth.1mg, hist(len[supp=="VC"],
12                        main="Histogram of length for Ascorbic Acid",
13                        xlab="Length"))
```

**Histogram of length for Orange Juice**



**Histogram of length for Ascorbic Acid**



```
1  # Formula interface
2  var.test(len ~ supp, data=Tooth.lmg, alternative="two.sided")
3  ##
4  ## F test to compare two variances
5  ##
6  ## data: len by supp
7  ## F = 2.4176, num df = 9, denom df = 9, p-value = 0.2046
8  ## alternative hypothesis: true ratio of variances is not equal
   to 1
9  ## 95 percent confidence interval:
10 ## 0.6004952 9.7332038
11 ## sample estimates:
12 ## ratio of variances
13 ##          2.41759
```



# 线性模型

## ◆ lm()

```
1 fit1 = lm(formula, data)
```

## ◆ 公式

## ❖ 公式构成

- 

变量名      公式所包含的变量名称，无需加引号

- 

波浪号 (~)      表示响应变量与控制变量的关系。

- 

加号 (+)      变量之间的线性关系。

- 

零 (0)      表示去除截距项。

- 

-      从等式中移除某个变量

•

• 数据集中除因变量以外的所有变量。

## ❖ 公式构成

- 

| 条件依赖。如  $y \sim \text{year} | \text{sex}$

- 

### Identity function

变量表达式。例如  $a+b$  表示变量  $a$  和  $b$  同时出现在公式中。 $I(a+b)$  表示公式中包含的是一个新变量  $= a + b$ 。

- 

: 自变量交互项

-

星号 (\*) 自变量所有可能交互项。例如： $y \sim a * b$  等价于  $y \sim a + b + I(a * b)$ ， $y \sim (a+b)*w$  等价于  $y \sim a + b + w + I(a*w) + I(b*w)$ 。

•  
 $\sim^n$  所有主效应及n阶交互项。例如  $y \sim (x + z + w)^2$  等价于  $y \sim x + z + w + x:z + x:w + z:w$

•  
 其它函数 例如：

1  $y \sim \log(u) + \sin(v) + w$

## ❖ 公式创建

```
1 sample.formula<-as.formula(y~x1+x2+x3)
2 class(sample.formula)
3 ## [1] "formula"
4 typeof(sample.formula)
5 ## [1] "language"
6
7 sample.formula2<-as.formula('y ~ x1 + x2 + x3')
8 class(sample.formula2)
9 ## [1] "formula"
10 typeof(sample.formula2)
11 ## [1] "language"
```



## ◆ 公式示例

Model	Interpretation
$y \sim 1$	仅有截距项
$y \sim a$	Just the intercept 一个主效应
$y \sim -1$	去除截距项
$+ a$	两个主效应
$y \sim a + b$	三个主效应，一个 a 和 b 的交互效应
$a+b+c+a:b$	a和b的主效应及其交互效应
$y \sim a*b$	
$\sim a + b + a:b$	
$y \sim \text{factor}(a)$	
$y \sim (a+b+c)^2$	
$y \sim I(a^2)$	
$\log(y) \sim y$	
$\sim .$	

## ◆ 简单线性回归

## ❖ 一元线性回归

```
1 # 萼片长度 ~ 萼片宽度
2 out <- lm(Sepal.Length ~ Sepal.Width, iris)
3 summary(out)
4 ##
5 ## Call:
6 ## lm(formula = Sepal.Length ~ Sepal.Width, data = iris)
7 ##
8 ## Residuals:
9 ##      Min       1Q   Median       3Q      Max
10 ## -1.5561 -0.6333 -0.1120  0.5579  2.2226
11 ##
12 ## Coefficients:
13 ##              Estimate Std. Error t value Pr(>|t|)
14 ## (Intercept)    6.5262     0.4789   13.63  <2e-16 ***
15 ## Sepal.Width   -0.2234     0.1551   -1.44    0.152
16 ## ---
```

```
17  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
18  ##
19  ## Residual standard error: 0.8251 on 148 degrees of freedom
20  ## Multiple R-squared:  0.01382,    Adjusted R-squared:  0.007159
21  ## F-statistic: 2.074 on 1 and 148 DF,  p-value: 0.1519
```

## ❖ 提取回归模型参数

```
1 summary(out)$r.squared
2 ## [1] 0.01382265
3 coeffs = coefficients(out); coeffs
4 ## (Intercept) Sepal.Width
5 ##      6.5262226  -0.2233611
```

```
1  # 方差协方差矩阵
2  vcov(out)
3  ##              (Intercept) Sepal.Width
4  ## (Intercept)  0.22934170 -0.07352916
5  ## Sepal.Width -0.07352916  0.02405009
6  # 估计系数的置信区间
7  confint(out)
8  ##              2.5 %      97.5 %
9  ## (Intercept)  5.579865 7.47258038
10 ## Sepal.Width -0.529820 0.08309785
```

## ❖ 估计因变量

```

1 fitted(out)
2 ##          1          2          3          4          5          6          7
3 ## 5.744459 5.856139 5.811467 5.833803 5.722123 5.655114 5.766795
4 ##          9          10         11         12         13         14         15
5 ## 5.878475 5.833803 5.699787 5.766795 5.856139 5.856139 5.632778
6 ##          17          18          19          20          21          22          23
7 ## 5.655114 5.744459 5.677451 5.677451 5.766795 5.699787 5.722123
8 ##          25          26          27          28          29          30          31
9 ##          32

```

```

9    ## 5.766795 5.856139 5.766795 5.744459 5.766795 5.811467 5.833803
    5.766795
10   ##          33          34          35          36          37          38          39
    40
11   ## 5.610442 5.588106 5.833803 5.811467 5.744459 5.722123 5.856139
    5.766795
12   ##          41          42          43          44          45          46          47
    48
13   ## 5.744459 6.012492 5.811467 5.744459 5.677451 5.856139 5.677451
    5.811467
14   ##          49          50          51          52          53          54          55
    56
15   ## 5.699787 5.789131 5.811467 5.811467 5.833803 6.012492 5.900812
    5.900812
16   ##          57          58          59          60          61          62          63
    64
17   ## 5.789131 5.990156 5.878475 5.923148 6.079500 5.856139 6.034828
    5.878475

```



18	##	65	66	67	68	69	70	71
	72							
19	##	5.878475	5.833803	5.856139	5.923148	6.034828	5.967820	5.811467
		5.900812						
20	##	73	74	75	76	77	78	79
	80							
21	##	5.967820	5.900812	5.878475	5.856139	5.900812	5.856139	5.878475
		5.945484						
22	##	81	82	83	84	85	86	87
	88							
23	##	5.990156	5.990156	5.923148	5.923148	5.856139	5.766795	5.833803
		6.012492						
24	##	89	90	91	92	93	94	95
	96							
25	##	5.856139	5.967820	5.945484	5.856139	5.945484	6.012492	5.923148
		5.856139						
26	##	97	98	99	100	101	102	103
	104							

```

27  ## 5.878475 5.878475 5.967820 5.900812 5.789131 5.923148 5.856139
    5.878475
28  ##      105      106      107      108      109      110      111
    112
29  ## 5.856139 5.856139 5.967820 5.878475 5.967820 5.722123 5.811467
    5.923148
30  ##      113      114      115      116      117      118      119
    120
31  ## 5.856139 5.967820 5.900812 5.811467 5.856139 5.677451 5.945484
    6.034828
32  ##      121      122      123      124      125      126      127
    128
33  ## 5.811467 5.900812 5.900812 5.923148 5.789131 5.811467 5.900812
    5.856139
34  ##      129      130      131      132      133      134      135
    136
35  ## 5.900812 5.856139 5.900812 5.677451 5.900812 5.900812 5.945484
    5.856139

```

```

36  ##      137      138      139      140      141      142      143
    144
37  ## 5.766795 5.833803 5.856139 5.833803 5.833803 5.833803 5.923148
    5.811467
38  ##      145      146      147      148      149      150
39  ## 5.789131 5.856139 5.967820 5.856139 5.766795 5.856139
40  ###predict(out)

```

## ❖ 估计因变量

```
1 # 置信区间估计
2 predict(out, iris, interval="confidence")
3 ##          fit          lwr          upr
4 ## 1  5.744459 5.554388 5.934529
5 ## 2  5.856139 5.721856 5.990423
6 ## 3  5.811467 5.671342 5.951592
7 ## 4  5.833803 5.700034 5.967573
8 ## 5  5.722123 5.509095 5.935150
9 ## 6  5.655114 5.364576 5.945653
10 ## 7  5.766795 5.597233 5.936357
11 ## 8  5.766795 5.597233 5.936357
12 ## 9  5.878475 5.736884 6.020067
13 ## 10 5.833803 5.700034 5.967573
14 ## 11 5.699787 5.462062 5.937511
15 ## 12 5.766795 5.597233 5.936357
16 ## 13 5.856139 5.721856 5.990423
```

17	##	14	5.856139	5.721856	5.990423
18	##	15	5.632778	5.314690	5.950866
19	##	16	5.543434	5.110961	5.975907
20	##	17	5.655114	5.364576	5.945653
21	##	18	5.744459	5.554388	5.934529
22	##	19	5.677451	5.413777	5.941124
23	##	20	5.677451	5.413777	5.941124
24	##	21	5.766795	5.597233	5.936357
25	##	22	5.699787	5.462062	5.937511
26	##	23	5.722123	5.509095	5.935150
27	##	24	5.789131	5.636639	5.941623
28	##	25	5.766795	5.597233	5.936357
29	##	26	5.856139	5.721856	5.990423
30	##	27	5.766795	5.597233	5.936357
31	##	28	5.744459	5.554388	5.934529
32	##	29	5.766795	5.597233	5.936357
33	##	30	5.811467	5.671342	5.951592
34	##	31	5.833803	5.700034	5.967573

35	##	32	5.766795	5.597233	5.936357
36	##	33	5.610442	5.264284	5.956601
37	##	34	5.588106	5.213473	5.962739
38	##	35	5.833803	5.700034	5.967573
39	##	36	5.811467	5.671342	5.951592
40	##	37	5.744459	5.554388	5.934529
41	##	38	5.722123	5.509095	5.935150
42	##	39	5.856139	5.721856	5.990423
43	##	40	5.766795	5.597233	5.936357
44	##	41	5.744459	5.554388	5.934529
45	##	42	6.012492	5.744929	6.280055
46	##	43	5.811467	5.671342	5.951592
47	##	44	5.744459	5.554388	5.934529
48	##	45	5.677451	5.413777	5.941124
49	##	46	5.856139	5.721856	5.990423
50	##	47	5.677451	5.413777	5.941124
51	##	48	5.811467	5.671342	5.951592
52	##	49	5.699787	5.462062	5.937511

53	##	50	5.789131	5.636639	5.941623
54	##	51	5.811467	5.671342	5.951592
55	##	52	5.811467	5.671342	5.951592
56	##	53	5.833803	5.700034	5.967573
57	##	54	6.012492	5.744929	6.280055
58	##	55	5.900812	5.746078	6.055546
59	##	56	5.900812	5.746078	6.055546
60	##	57	5.789131	5.636639	5.941623
61	##	58	5.990156	5.748694	6.231618
62	##	59	5.878475	5.736884	6.020067
63	##	60	5.923148	5.750766	6.095529
64	##	61	6.079500	5.729189	6.429812
65	##	62	5.856139	5.721856	5.990423
66	##	63	6.034828	5.740287	6.329369
67	##	64	5.878475	5.736884	6.020067
68	##	65	5.878475	5.736884	6.020067
69	##	66	5.833803	5.700034	5.967573
70	##	67	5.856139	5.721856	5.990423

71	##	68	5.923148	5.750766	6.095529
72	##	69	6.034828	5.740287	6.329369
73	##	70	5.967820	5.751265	6.184375
74	##	71	5.811467	5.671342	5.951592
75	##	72	5.900812	5.746078	6.055546
76	##	73	5.967820	5.751265	6.184375
77	##	74	5.900812	5.746078	6.055546
78	##	75	5.878475	5.736884	6.020067
79	##	76	5.856139	5.721856	5.990423
80	##	77	5.900812	5.746078	6.055546
81	##	78	5.856139	5.721856	5.990423
82	##	79	5.878475	5.736884	6.020067
83	##	80	5.945484	5.752180	6.138788
84	##	81	5.990156	5.748694	6.231618
85	##	82	5.990156	5.748694	6.231618
86	##	83	5.923148	5.750766	6.095529
87	##	84	5.923148	5.750766	6.095529
88	##	85	5.856139	5.721856	5.990423



89	##	86	5.766795	5.597233	5.936357
90	##	87	5.833803	5.700034	5.967573
91	##	88	6.012492	5.744929	6.280055
92	##	89	5.856139	5.721856	5.990423
93	##	90	5.967820	5.751265	6.184375
94	##	91	5.945484	5.752180	6.138788
95	##	92	5.856139	5.721856	5.990423
96	##	93	5.945484	5.752180	6.138788
97	##	94	6.012492	5.744929	6.280055
98	##	95	5.923148	5.750766	6.095529
99	##	96	5.856139	5.721856	5.990423
100	##	97	5.878475	5.736884	6.020067
101	##	98	5.878475	5.736884	6.020067
102	##	99	5.967820	5.751265	6.184375
103	##	100	5.900812	5.746078	6.055546
104	##	101	5.789131	5.636639	5.941623
105	##	102	5.923148	5.750766	6.095529
106	##	103	5.856139	5.721856	5.990423

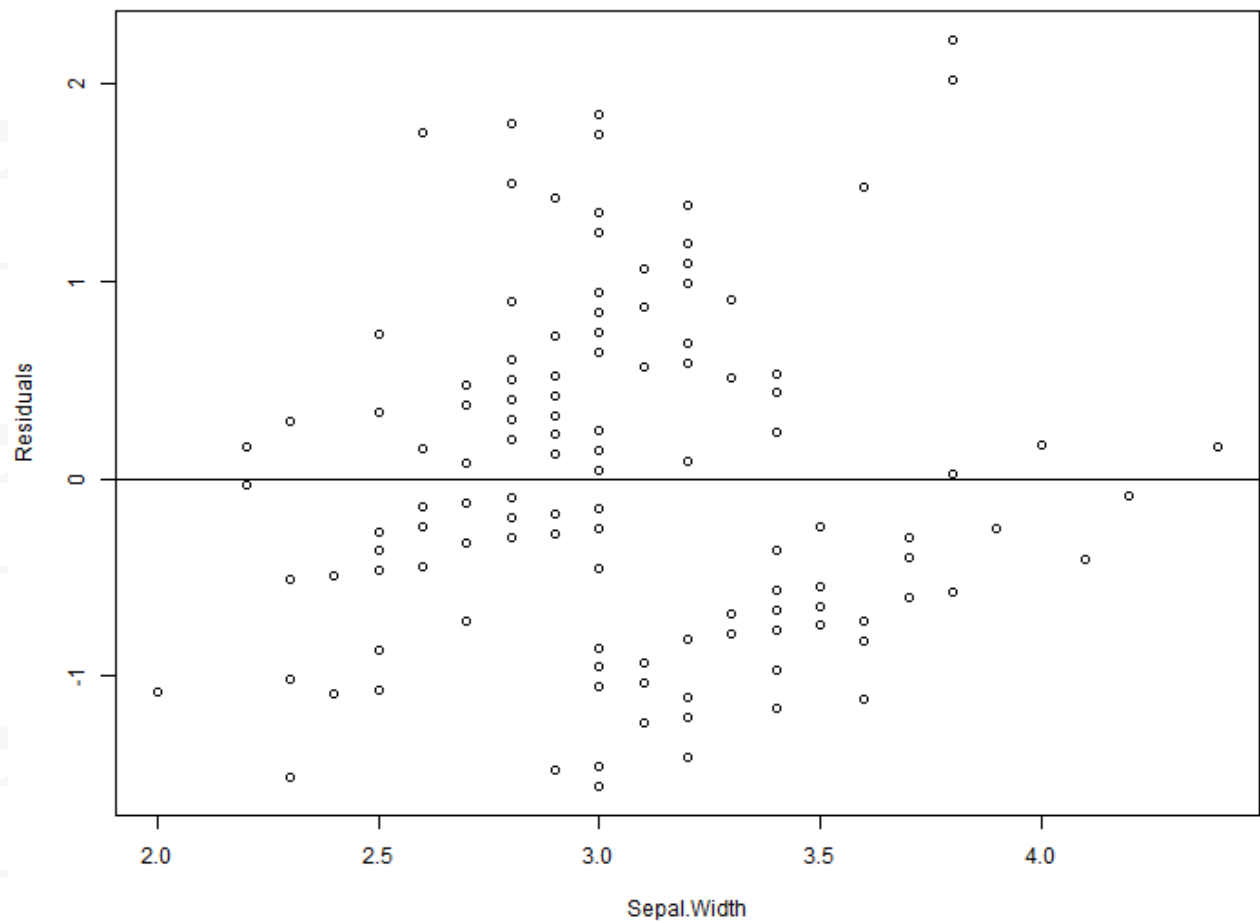
```
107  ## 104 5.878475 5.736884 6.020067
108  ## 105 5.856139 5.721856 5.990423
109  ## 106 5.856139 5.721856 5.990423
110  ## 107 5.967820 5.751265 6.184375
111  ## 108 5.878475 5.736884 6.020067
112  ## 109 5.967820 5.751265 6.184375
113  ## 110 5.722123 5.509095 5.935150
114  ## 111 5.811467 5.671342 5.951592
115  ## 112 5.923148 5.750766 6.095529
116  ## 113 5.856139 5.721856 5.990423
117  ## 114 5.967820 5.751265 6.184375
118  ## 115 5.900812 5.746078 6.055546
119  ## 116 5.811467 5.671342 5.951592
120  ## 117 5.856139 5.721856 5.990423
121  ## 118 5.677451 5.413777 5.941124
122  ## 119 5.945484 5.752180 6.138788
123  ## 120 6.034828 5.740287 6.329369
124  ## 121 5.811467 5.671342 5.951592
```

125	##	122	5.900812	5.746078	6.055546
126	##	123	5.900812	5.746078	6.055546
127	##	124	5.923148	5.750766	6.095529
128	##	125	5.789131	5.636639	5.941623
129	##	126	5.811467	5.671342	5.951592
130	##	127	5.900812	5.746078	6.055546
131	##	128	5.856139	5.721856	5.990423
132	##	129	5.900812	5.746078	6.055546
133	##	130	5.856139	5.721856	5.990423
134	##	131	5.900812	5.746078	6.055546
135	##	132	5.677451	5.413777	5.941124
136	##	133	5.900812	5.746078	6.055546
137	##	134	5.900812	5.746078	6.055546
138	##	135	5.945484	5.752180	6.138788
139	##	136	5.856139	5.721856	5.990423
140	##	137	5.766795	5.597233	5.936357
141	##	138	5.833803	5.700034	5.967573
142	##	139	5.856139	5.721856	5.990423

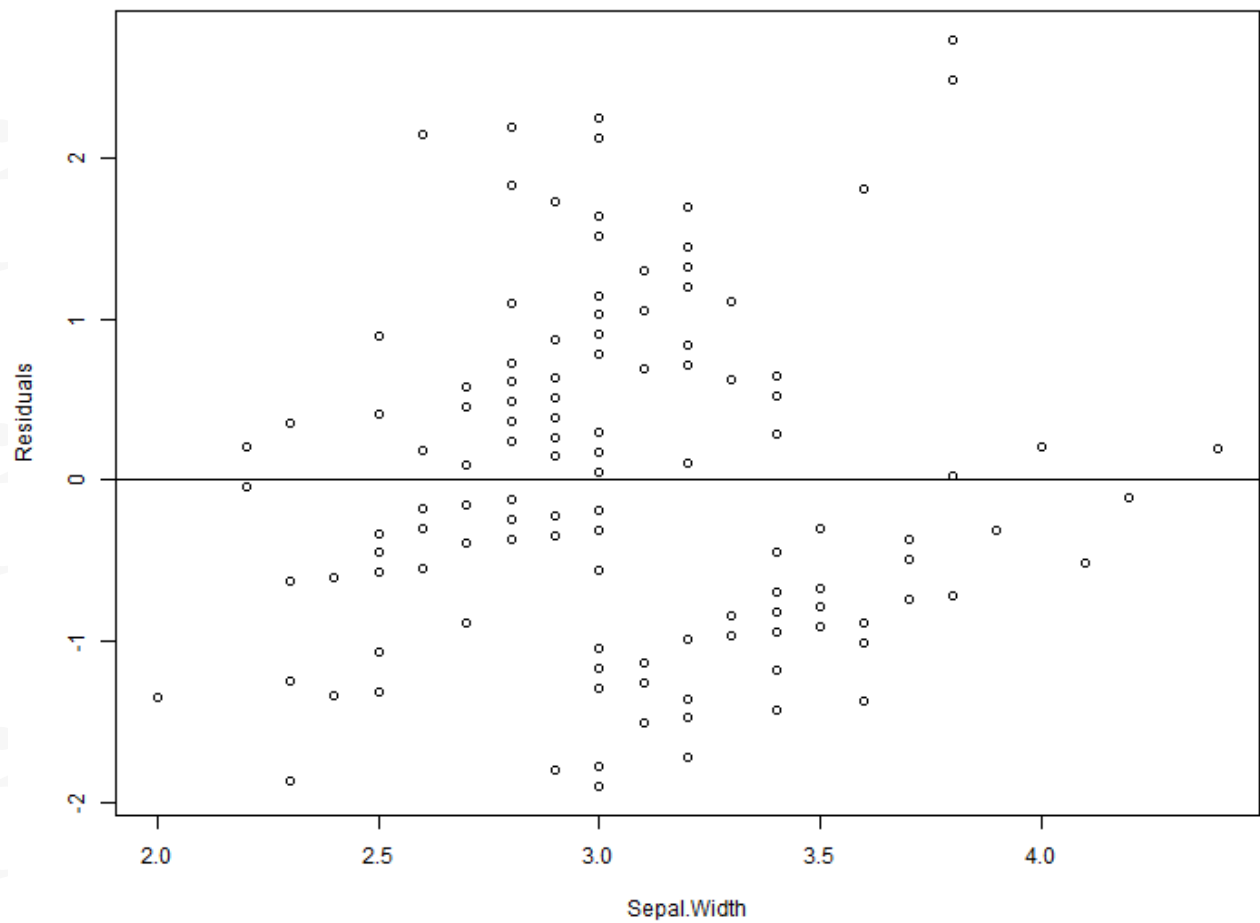
```
143    ## 140 5.833803 5.700034 5.967573
144    ## 141 5.833803 5.700034 5.967573
145    ## 142 5.833803 5.700034 5.967573
146    ## 143 5.923148 5.750766 6.095529
147    ## 144 5.811467 5.671342 5.951592
148    ## 145 5.789131 5.636639 5.941623
149    ## 146 5.856139 5.721856 5.990423
150    ## 147 5.967820 5.751265 6.184375
151    ## 148 5.856139 5.721856 5.990423
152    ## 149 5.766795 5.597233 5.936357
153    ## 150 5.856139 5.721856 5.990423
```

## ❖ 残差

```
1 lm.res = resid(out)
2 lm.res2 = iris$Sepal.Length - fitted(out)
3 all.equal(lm.res, lm.res2, check.attributes = F)
4 ## [1] TRUE
5
6 plot(iris$Sepal.Width, lm.res,
7      ylab="Residuals", xlab="Sepal.Width",
8      main="")
9 abline(0, 0) # the horizon
```



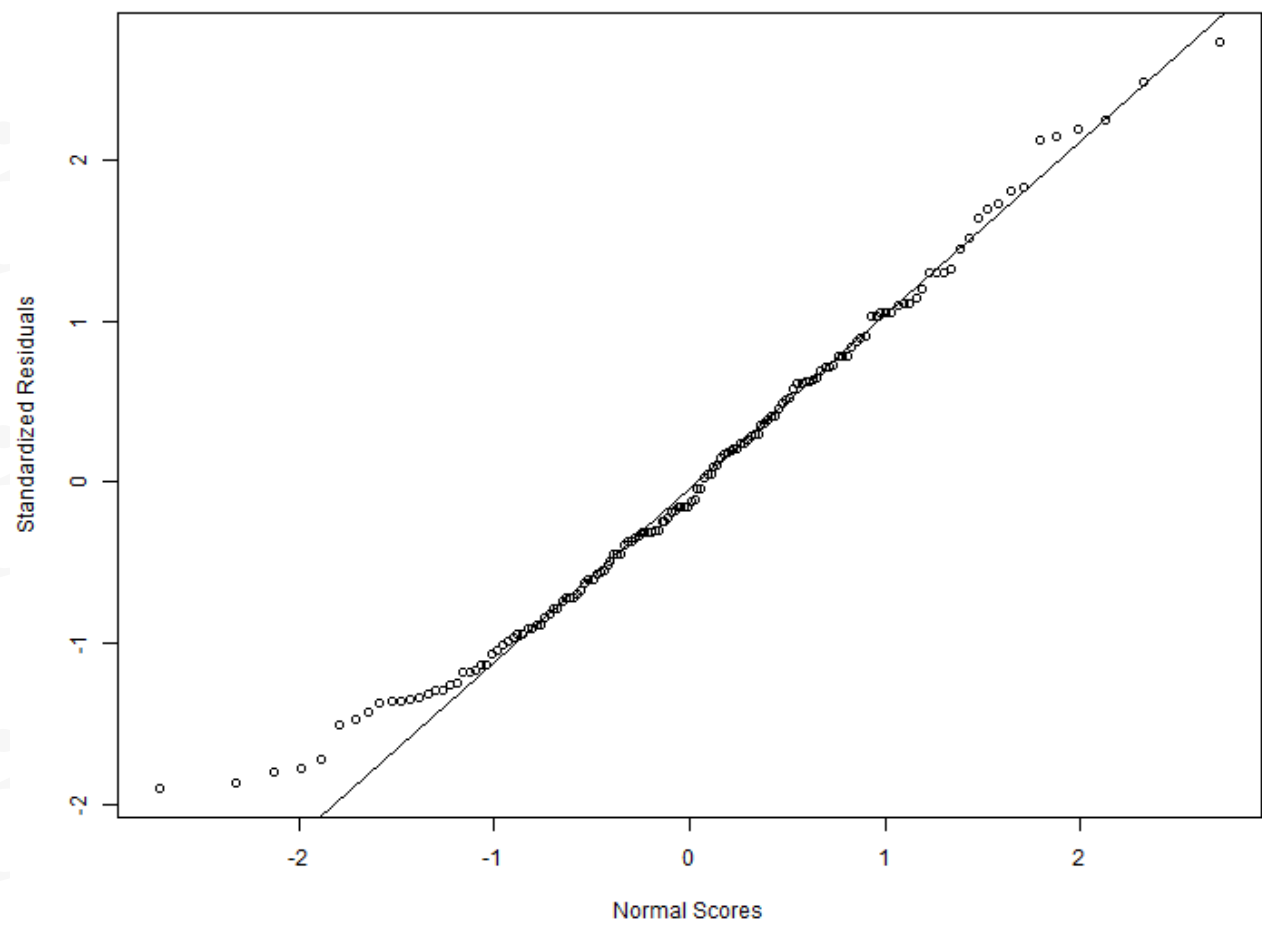
```
1 # 标准化残差 $res.std = res/ std(res)$  
2 stdres<-rstandard(out)  
3 plot(iris$Sepal.Width, stdres,  
4      ylab="Residuals", xlab="Sepal.Width",  
5      main="")  
6 abline(0, 0)
```





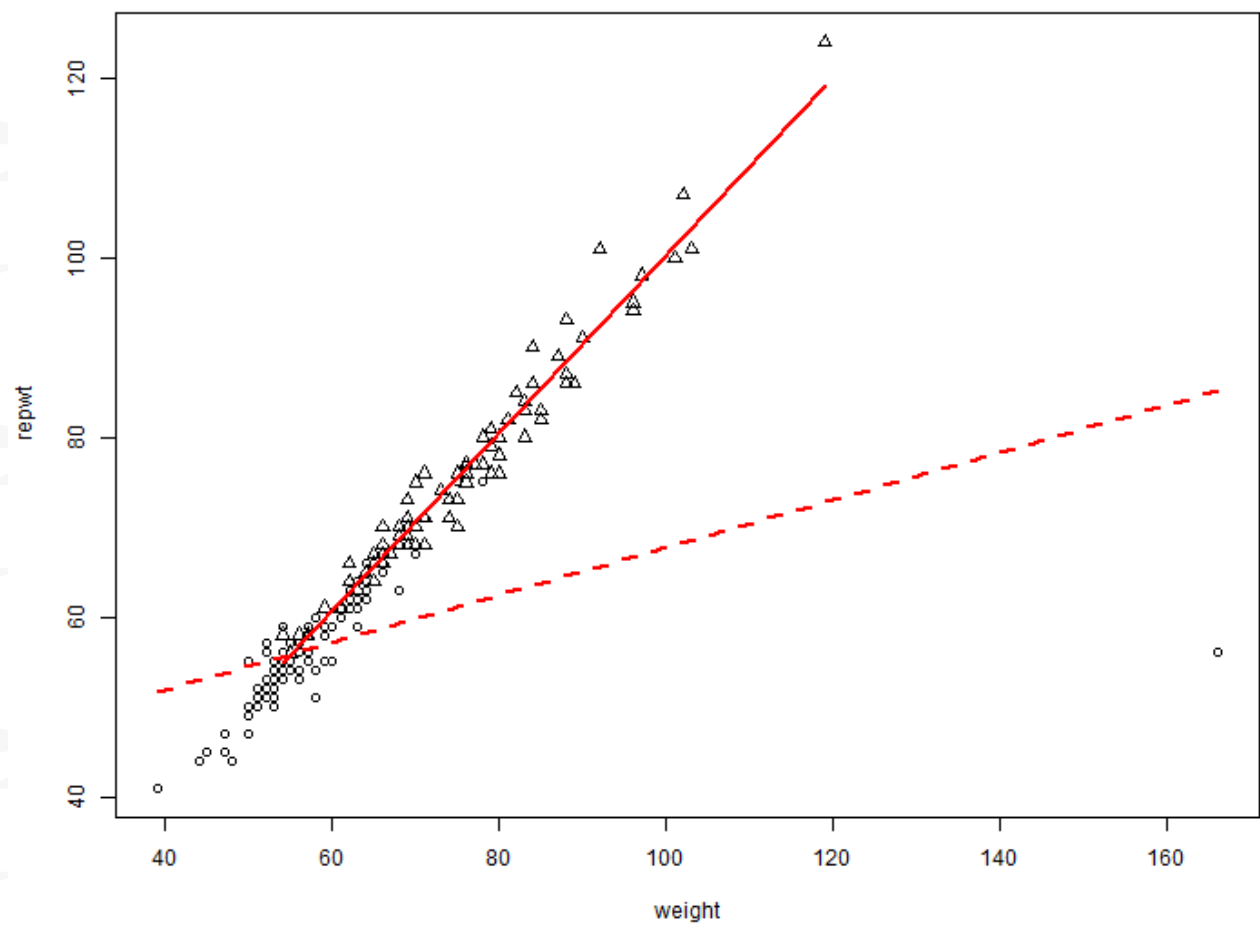
## ❖ 残差正态概率图

```
1 qqnorm(stdres,  
2       ylab="Standardized Residuals",  
3       xlab="Normal Scores",  
4       main="")  
5 qqline(stdres)
```



## ❖ 使用 car 包制图

```
1 library(car)
2 plot(repwt ~ weight, pch=c(1,2)[sex], data=Davis)
3 regLine(lm(repwt~weight, subset=sex=="M", data=Davis))
4 regLine(lm(repwt~weight, subset=sex=="F", data=Davis), lty=2)
```



## ❖ 异常值

```
1 plot(iris$Sepal.Width, iris$Sepal.Length)
2 identify(iris$Sepal.Width, iris$Sepal.Length)
```

进入交互模式，用鼠标左键逐个点击散点图上的异常值，完成后按 **Esc** 或鼠标右键退出。**R** 自动输出被标识为异常值的序号。

## ❖ 标准化回归系数

- 方式一: 对数据进行标准化处理, 将数据转化为均值为0, 标准差为1, 然后再用 `lm()` 函数拟合。

```
1 iris_std = as.data.frame(scale(iris[, c('Sepal.Length', 'Sepal.Width')]))
2 out_std <- lm(Sepal.Length ~ Sepal.Width, iris_std)
3 summary(out_std)
4 ##
5 ## Call:
6 ## lm(formula = Sepal.Length ~ Sepal.Width, data = iris_std)
7 ##
8 ## Residuals:
9 ##      Min       1Q   Median       3Q      Max
10 ## -1.8793 -0.7648 -0.1352  0.6738  2.6840
11 ##
12 ## Coefficients:
13 ##              Estimate Std. Error t value Pr(>|t|)
```

```
14  ## (Intercept) -3.759e-16  8.136e-02    0.00    1.000
15  ## Sepal.Width -1.176e-01  8.163e-02   -1.44    0.152
16  ##
17  ## Residual standard error: 0.9964 on 148 degrees of freedom
18  ## Multiple R-squared:  0.01382,    Adjusted R-squared:  0.007159
19  ## F-statistic: 2.074 on 1 and 148 DF,  p-value: 0.1519
```

## ❖ 标准化回归系数

- 方式二：在设定回归方程的时候，去掉截距项。

由于少用了一个自由度，除了自变量的回归系数以外，模型的残标准误、回归系数的标准误、t值都会发生变化。

```
1  lm(y ~ x - 1)
1  out3 <- lm(Sepal.Length ~ Sepal.Width -1, iris)
2  summary(out3)
3  ##
4  ## Call:
5  ## lm(formula = Sepal.Length ~ Sepal.Width - 1, data = iris)
6  ##
7  ## Residuals:
8  ##      Min       1Q   Median       3Q      Max
9  ## -2.5236 -1.0362  0.4823  0.9897  2.8406
10 ##
```

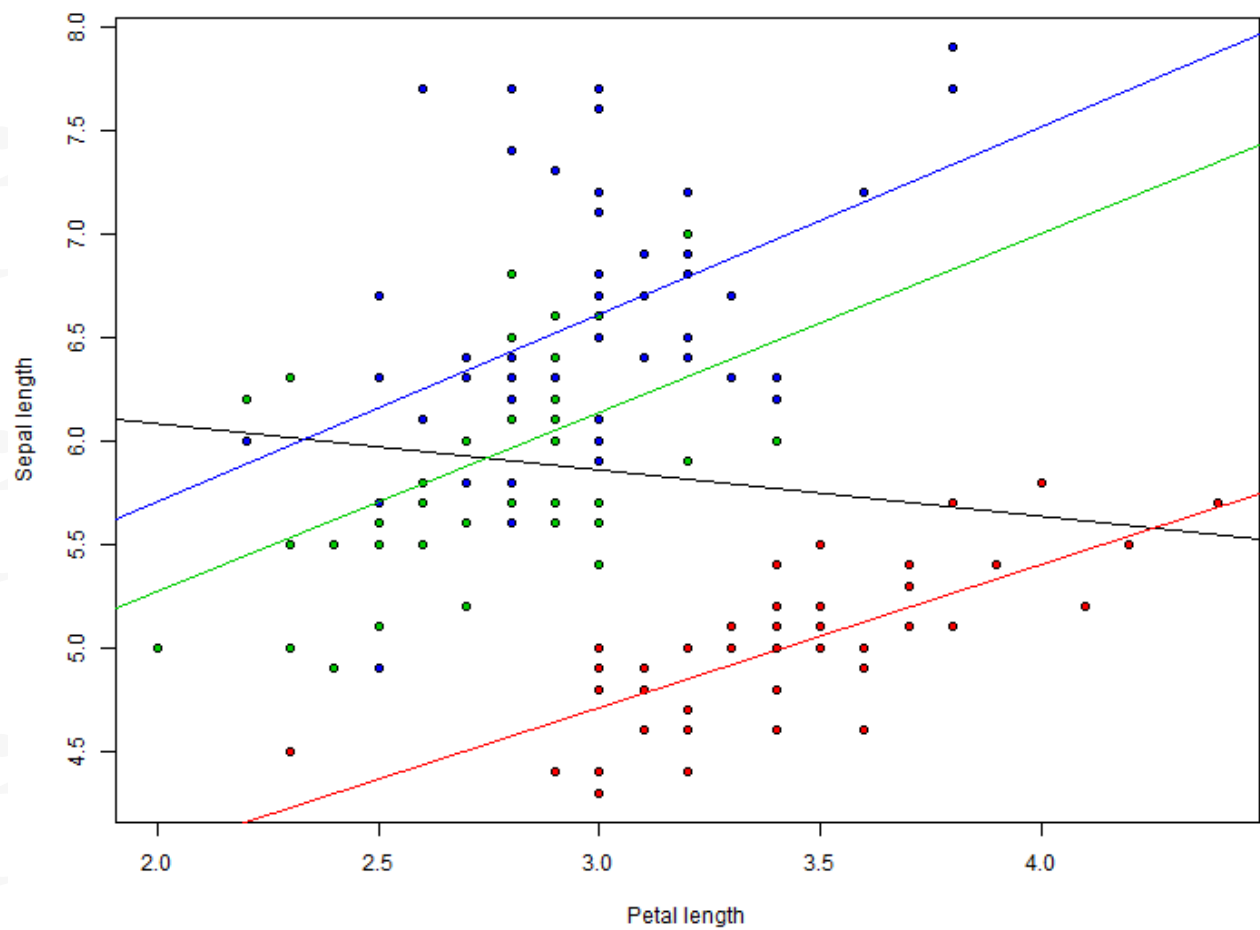


```
11  ## Coefficients:
12  ##           Estimate Std. Error t value Pr(>|t|)
13  ## Sepal.Width  1.86901    0.03265   57.25  <2e-16 ***
14  ## ---
15  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16  ##
17  ## Residual standard error: 1.235 on 149 degrees of freedom
18  ## Multiple R-squared:  0.9565, Adjusted R-squared:  0.9562
19  ## F-statistic: 3277 on 1 and 149 DF,  p-value: < 2.2e-16
```

以 `iris` 数据为例，做4个回归模型：第一个模型利用全部数据，另外三个分别用筛选后的数据。

```
1 plot(iris$Sepal.Width, iris$Sepal.Length, pch=21, bg=c("red","green3","blue"),
2      main="Edgar Anderson's Iris Data", xlab="Petal length", ylab="Sepal
3      length")
4 abline(lm(Sepal.Length ~ Sepal.Width, data=iris)$coefficients,
5         col="black")
6 abline(lm(Sepal.Length ~ Sepal.Width, data=iris[which(iris$Species=="setosa"),
7               ])$coefficients, col="red")
8 abline(lm(Sepal.Length ~ Sepal.Width, data=iris[which(iris$Species=="versicolour"),
9               ])$coefficients, col="green3")
10 abline(lm(Sepal.Length ~ Sepal.Width, data=iris[which(iris$Species=="virginica"),
11               ])$coefficients, col="blue")
```

Edgar Anderson's Iris Data



## ◆ 多元线性回归

## ❖ 示例

```
1  fit <- lm(len~factor(dose)+supp, data=ToothGrowth)
2  summary(fit)
3  ##
4  ## Call:
5  ## lm(formula = len ~ factor(dose) + supp, data = ToothGrowth)
6  ##
7  ## Residuals:
8  ##      Min       1Q   Median       3Q      Max
9  ## -7.085 -2.751 -0.800  2.446  9.650
10 ##
11 ## Coefficients:
12 ##              Estimate Std. Error t value Pr(>|t|)
13 ## (Intercept)    12.4550     0.9883  12.603 < 2e-16 ***
14 ## factor(dose)1     9.1300     1.2104   7.543 4.38e-10 ***
15 ## factor(dose)2    15.4950     1.2104  12.802 < 2e-16 ***
16 ## suppVC          -3.7000     0.9883  -3.744 0.000429 ***
```

```
17  ## ---
18  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
19  ##
20  ## Residual standard error: 3.828 on 56 degrees of freedom
21  ## Multiple R-squared:  0.7623, Adjusted R-squared:  0.7496
22  ## F-statistic: 59.88 on 3 and 56 DF,  p-value: < 2.2e-16
```

```

1  anova(fit)
2  ## Analysis of Variance Table
3  ##
4  ## Response: len
5  ##           Df   Sum Sq Mean Sq F value    Pr(>F)
6  ## factor(dose)  2 2426.43 1213.22   82.811 < 2.2e-16 ***
7  ## supp         1  205.35   205.35   14.017 0.0004293 ***
8  ## Residuals    56  820.43    14.65
9  ## ---
10 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

1 summary(fit)$r.squared
2 ## [1] 0.7623478
3 # $r.adj = 1- (1-R^2) (n-1)/(n-p-1)$
4 summary(fit)$adj.r.squared
5 ## [1] 0.7496165

1 ssq = anova(fit)$`Sum Sq`
2 sum(head(ssq,2))/sum(ssq)
3 ## [1] 0.7623478

1 dose.2 <- relevel(factor(ToothGrowth$dose), 3)
2 fit <- lm(len~dose.2+supp, data=ToothGrowth)
3 summary(fit)
4 ##
5 ## Call:
6 ## lm(formula = len ~ dose.2 + supp, data = ToothGrowth)
7 ##
8 ## Residuals:
9 ##      Min       1Q   Median       3Q      Max

```



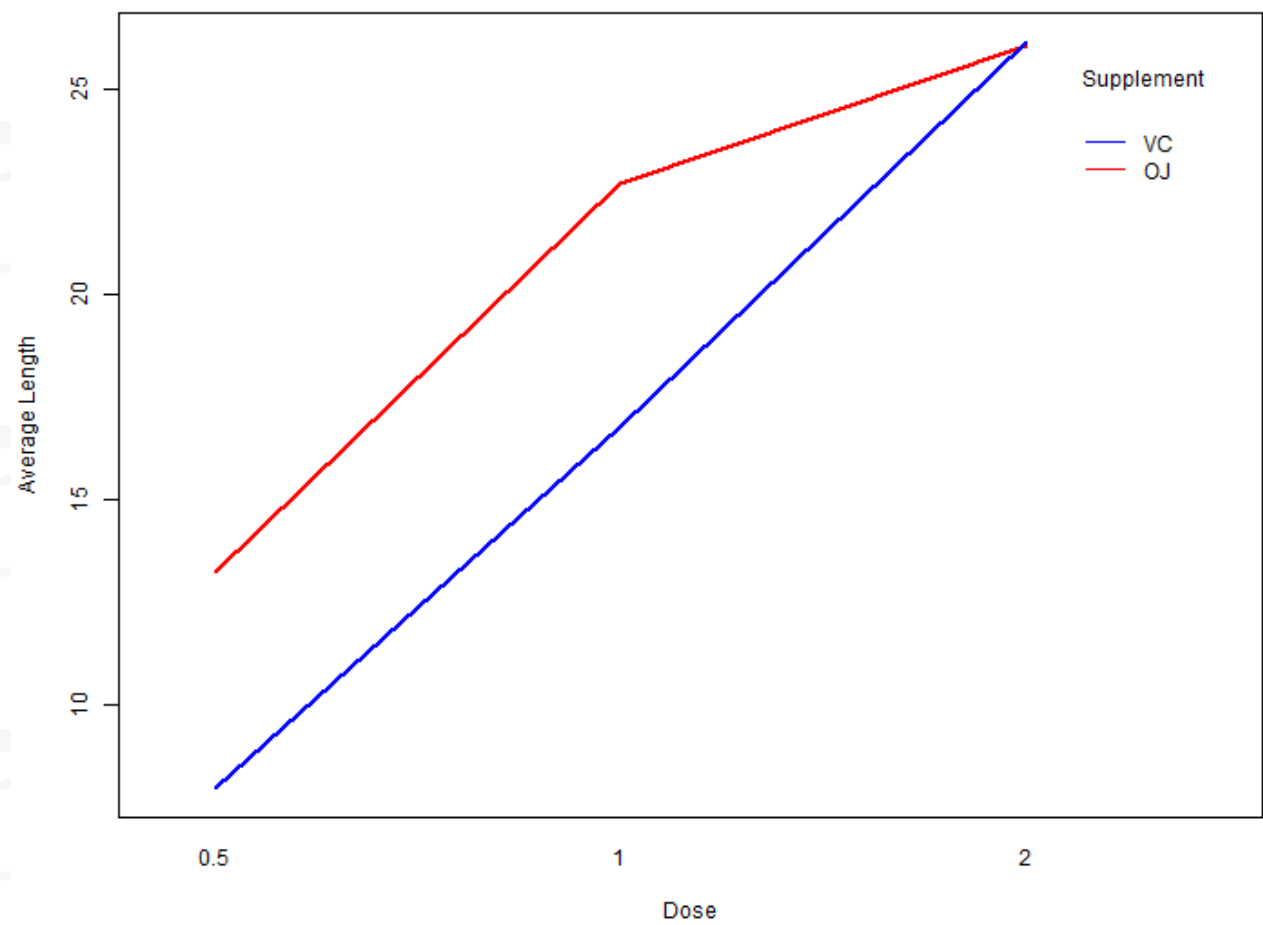
```

10  ## -7.085 -2.751 -0.800  2.446  9.650
11  ##
12  ## Coefficients:
13  ##           Estimate Std. Error t value Pr(>|t|)
14  ## (Intercept)  27.9500      0.9883  28.281 < 2e-16 ***
15  ## dose.20.5    -15.4950      1.2104 -12.802 < 2e-16 ***
16  ## dose.21      -6.3650      1.2104  -5.259 2.35e-06 ***
17  ## suppVC       -3.7000      0.9883  -3.744 0.000429 ***
18  ## ---
19  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
20  ##
21  ## Residual standard error: 3.828 on 56 degrees of freedom
22  ## Multiple R-squared:  0.7623, Adjusted R-squared:  0.7496
23  ## F-statistic: 59.88 on 3 and 56 DF,  p-value: < 2.2e-16

```

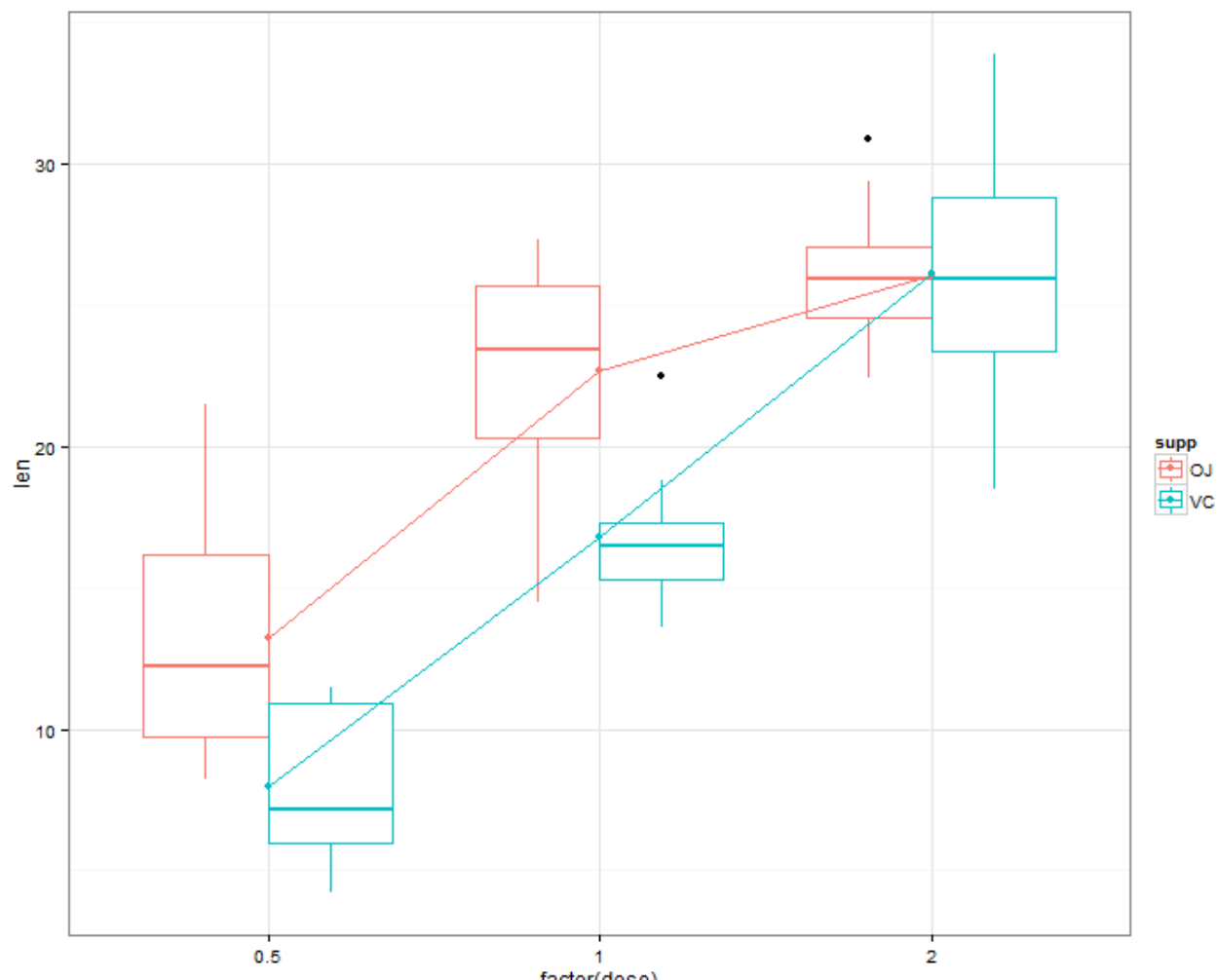
## ❖ 画图

```
1   with(ToothGrowth,  
2       interaction.plot(x.factor=dose, trace.factor=supp, response=len,  
3                       fun=mean,  
4                       xlab="Dose", ylab="Average Length", trace.label="Su  
5                       lty=1, lwd=2, col=c("red", "blue")))
```



## ❖ ggplot 实现

```
1   toothInt = aggregate(len ~ dose + supp, ToothGrowth, mean)
2
3   library(ggplot2)
4   ggplot(ToothGrowth, aes(x = factor(dose), y = len, colour = supp))
5   +
6     geom_boxplot() +
7     geom_point(data = toothInt, aes(y = len)) +
8     geom_line(data = toothInt, aes(y = len, group = supp)) +
9     theme_bw()
```



## ❖ 嵌套模型检验

什么是嵌套模型?

```
1  fit.1 = lm(len ~ factor(dose), ToothGrowth)
2  fit.2 = lm(len ~ factor(dose) + supp, ToothGrowth)
3  anova(fit.2, fit.1)
4  ## Analysis of Variance Table
5  ##
6  ## Model 1: len ~ factor(dose) + supp
7  ## Model 2: len ~ factor(dose)
8  ##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
9  ## 1      56  820.43
10 ## 2      57 1025.78 -1    -205.35 14.017 0.0004293 ***
11 ## ---
12 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

1 drop1(fit.2, ~supp, test="F")
2 ## Single term deletions
3 ##
4 ## Model:
5 ## len ~ factor(dose) + supp
6 ##      Df Sum of Sq    RSS    AIC F value    Pr(>F)
7 ## <none>          820.43 164.93
8 ## supp      1      205.35 1025.78 176.33   14.017 0.0004293 ***
9 ## ---
10 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

1 library(lmtest)
2 lrtest(m1, m2)

```

## ❖ 变量选择

```
1 step(object, scope, scale = 0, direction = c("both", "backward",
2   "forward"),
3   trace = 1, keep = NULL, steps = 1000, k = 2, ...)
1 summary(lm1 <- lm(Fertility ~ ., data = swiss))
2 ##
3 ## Call:
4 ## lm(formula = Fertility ~ ., data = swiss)
5 ##
6 ## Residuals:
7 ##      Min       1Q   Median       3Q      Max
8 ## -15.2743  -5.2617   0.5032   4.1198  15.3213
9 ##
10 ## Coefficients:
11 ##              Estimate Std. Error t value Pr(>|t|)
12 ## (Intercept)    66.91518    10.70604     6.250 1.91e-07 ***
13 ## Agriculture   -0.17211     0.07030    -2.448 0.01873 *
```



```

14  ## Examination      -0.25801    0.25388   -1.016    0.31546
15  ## Education        -0.87094    0.18303   -4.758 2.43e-05 ***
16  ## Catholic         0.10412    0.03526    2.953 0.00519 **
17  ## Infant.Mortality 1.07705    0.38172    2.822 0.00734 **
18  ## ---
19  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
20  ##
21  ## Residual standard error: 7.165 on 41 degrees of freedom
22  ## Multiple R-squared:  0.7067, Adjusted R-squared:  0.671
23  ## F-statistic: 19.76 on 5 and 41 DF,  p-value: 5.594e-10
24  slm1 <- step(lm1)
25  ## Start:  AIC=190.69
26  ## Fertility ~ Agriculture + Examination + Education + Catholic +
27  ##      Infant.Mortality
28  ##
29  ##              Df Sum of Sq    RSS    AIC
30  ## - Examination    1      53.03 2158.1 189.86
31  ## <none>              2105.0 190.69

```

```

32  ## - Agriculture      1    307.72 2412.8 195.10
33  ## - Infant.Mortality 1    408.75 2513.8 197.03
34  ## - Catholic        1    447.71 2552.8 197.75
35  ## - Education       1   1162.56 3267.6 209.36
36  ##
37  ## Step:  AIC=189.86
38  ## Fertility ~ Agriculture + Education + Catholic + Infant.Mortality
39  ##
40  ##               Df Sum of Sq    RSS    AIC
41  ## <none>                2158.1 189.86
42  ## - Agriculture      1    264.18 2422.2 193.29
43  ## - Infant.Mortality 1    409.81 2567.9 196.03
44  ## - Catholic        1    956.57 3114.6 205.10
45  ## - Education       1   2249.97 4408.0 221.43
46  summary(slm1)
47  ##
48  ## Call:
49  ## lm(formula = Fertility ~ Agriculture + Education + Catholic +

```

```

50  ##      Infant.Mortality, data = swiss)
51  ##
52  ## Residuals:
53  ##      Min      1Q    Median      3Q      Max
54  ## -14.6765  -6.0522   0.7514   3.1664  16.1422
55  ##
56  ## Coefficients:
57  ##              Estimate Std. Error t value Pr(>|t|)
58  ## (Intercept)    62.10131     9.60489   6.466 8.49e-08 ***
59  ## Agriculture    -0.15462     0.06819  -2.267  0.02857 *
60  ## Education      -0.98026     0.14814  -6.617 5.14e-08 ***
61  ## Catholic        0.12467     0.02889   4.315 9.50e-05 ***
62  ## Infant.Mortality 1.07844     0.38187   2.824 0.00722 **
63  ## ---
64  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
65  ##
66  ## Residual standard error: 7.168 on 42 degrees of freedom
67  ## Multiple R-squared:  0.6993, Adjusted R-squared:  0.6707

```

```

68  ## F-statistic: 24.42 on 4 and 42 DF,  p-value: 1.717e-10
69  slm1$anova
70  ##              Step Df Deviance Resid. Df Resid. Dev      AIC
71  ## 1              NA      NA      41    2105.043 190.6913
72  ## 2 - Examination  1  53.02656      42    2158.069 189.8606

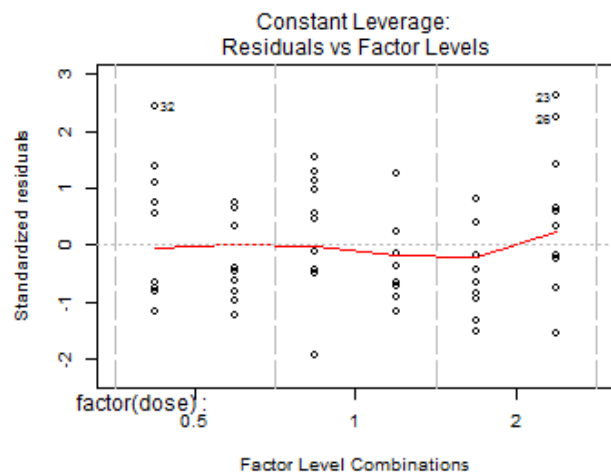
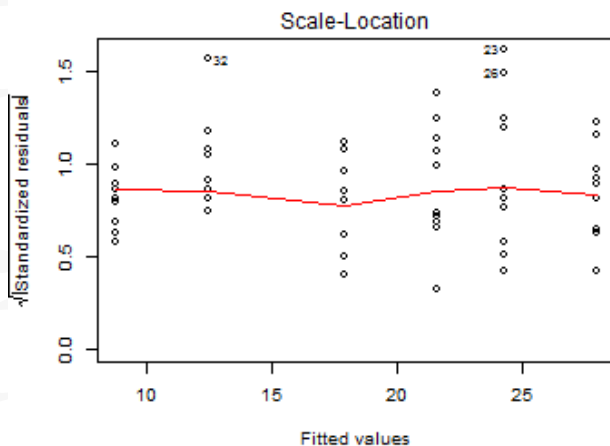
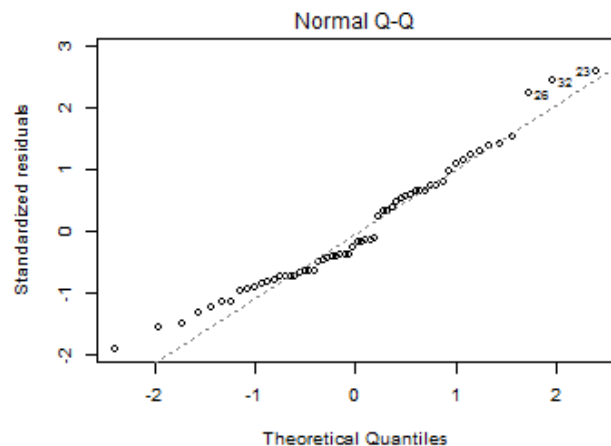
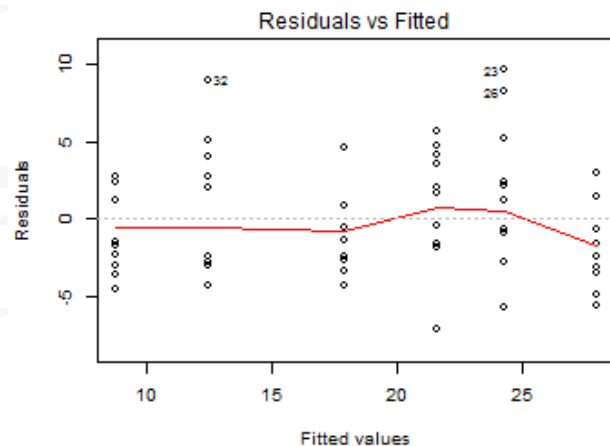
```

## ❖ 模型诊断

诊断函数	说明
<code>fitted.values()</code> <code>fitted( )</code> <code>residuals( )</code> <code>rstandard( )</code> <code>rstudent( )</code> <code>qqnorm()</code> <code>qqline()</code> <code>plot.lm()</code>	

## ❖ 模型诊断示例

```
1 fit <- lm(len~factor(dose)+supp, data=ToothGrowth)
2 par(mfrow=c(2,2))
3 plot(fit)
```



## ❖ 影响点判断

```
1      # dfbetas, dffits, covratio and cooks.distance
2      inf.temp <- influence.measures(fit)
3      inf.pts <- which(apply(inf.temp$is.inf, 1, any))
4      ToothGrowth[inf.pts,]
5      ##      len supp dose
6      ## 23 33.9   VC   2.0
7      ## 32 21.5   OJ   0.5
```



```

1  lm.inf.coef <- lm.influence(fit)$coefficients
2  lm.inf.pts <- apply(lm.inf.coef, 2, FUN=function(x) which.max(abs(x)))
3  # Get the five points that cause the greatest change in the estimates
4  lm.inf.pts.top5 <- apply(lm.inf.coef, 2, FUN=function(x)
5    names(rev(sort(abs(x)))[1:5]))
6  lm.inf.pts.top5
7  ##      (Intercept) factor(dose)1 factor(dose)2 suppVC
8  ## [1,] "32"          "32"          "23"          "23"
9  ## [2,] "33"          "49"          "32"          "32"
10 ## [3,] "37"          "50"          "26"          "26"
11 ## [4,] "39"          "33"          "22"          "49"
12 ## [5,] "38"          "44"          "53"          "22"

```

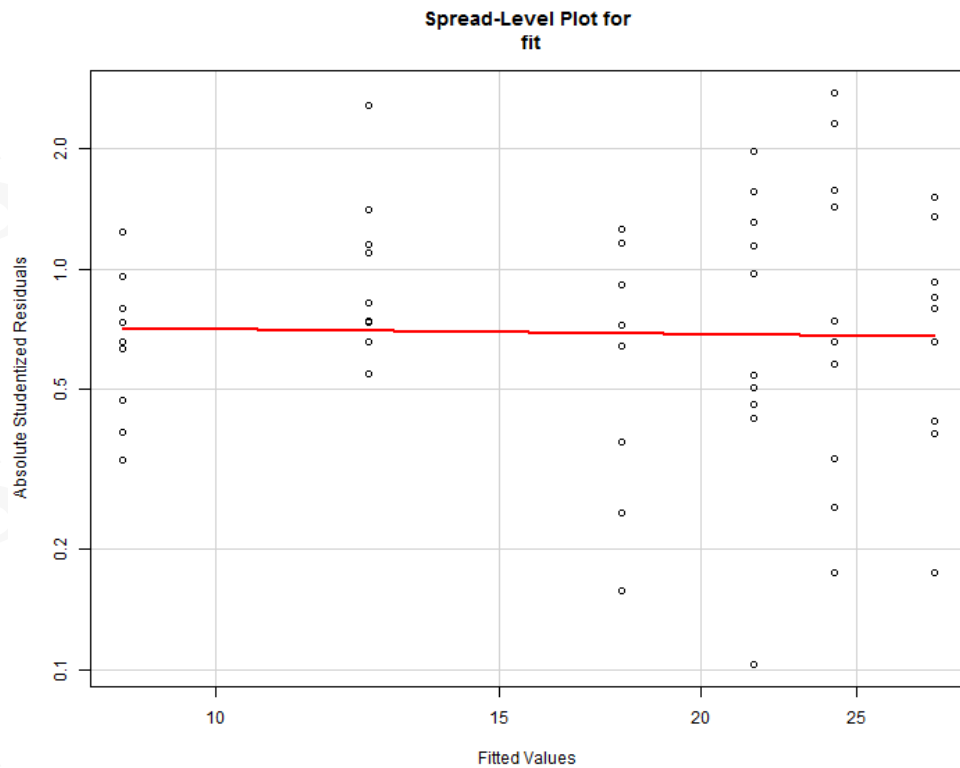
## ❖ 方差同质性检验

- car:

- `ncv.test()`: 生成一个计分检验, 零假设为误差方差守恒, 备择假设则为误差方差随拟合值变化而变化。
- `spreadLevelPlot()`: 包含最佳拟合曲线的散点图。

```
1 library(car)
2 ncvTest(fit)
3 ## Non-constant Variance Score Test
4 ## Variance formula: ~ fitted.values
5 ## Chisquare = 0.4405518    Df = 1    p = 0.5068562
```

1 spreadLevelPlot(fit)



```
1      ##  
2      ## Suggested power transformation:  1.037044
```

## ◆ 广义线性模型

- Generalized Linear Models are fit using the function `glm()`.  
`glm(formula, family = gaussian, data)`
- `family` 参数设定分布类型和链接函数

```
1  binomial(link = "logit")
2  ##
3  ## Family: binomial
4  ## Link function: logit
5  gaussian(link = "identity")
6  ##
7  ## Family: gaussian
8  ## Link function: identity
9  poisson(link = "log")
10 ##
11 ## Family: poisson
12 ## Link function: log
```

---

<code>summary.glm()</code>	Summarize the model fit
<code>anova.glm()</code>	Analysis of deviance table
<code>confint.glm()</code>	Confidence interval for model parameters
<code>predict.glm()</code>	Obtain predicted values
<code>influence.measures()</code>	Measures of influence
<code>step()</code>	Step-wise selection using AIC
<code>drop1()</code>	Test parameter using deviance

---

## ❖ 示例

```
1 y <- ifelse(ToothGrowth[,1]>20, 1, 0)
2 # Fit logistic model
3 fit <- glm(y~supp+factor(dose), family="binomial", data=ToothGrowth)
4 summary(fit)
5 ##
6 ## Call:
7 ## glm(formula = y ~ supp + factor(dose)      , family = "binomial",
8 data = ToothGrowth)
9 ##
10 ## Deviance Residuals:
11 ##      Min       1Q   Median       3Q      Max
12 ## -2.16659  -0.43759  -0.09337   0.44842   2.18796
13 ##
14 ## Coefficients:
15 ##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -2.247       1.051  -2.138  0.03253 *
```

```

16  ## suppVC          -3.187      1.218  -2.617  0.00887 **
17  ## factor(dose)1    3.136      1.237   2.534  0.01128 *
18  ## factor(dose)2    7.680      1.848   4.156  3.24e-05 ***
19  ## ---
20  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
21  ##
22  ## (Dispersion parameter for binomial family taken to be 1)
23  ##
24  ##      Null deviance: 82.911  on 59  degrees of freedom
25  ## Residual deviance: 31.913  on 56  degrees of freedom
26  ## AIC: 39.913
27  ##
28  ## Number of Fisher Scoring iterations: 6

```



```

1  anova(fit, test="Chisq") # Compare reduction in deviance, sequentially
2  ## Analysis of Deviance Table
3  ##
4  ## Model: binomial, link: logit
5  ##
6  ## Response: y
7  ##
8  ## Terms added sequentially (first to last)
9  ##
10 ##
11 ##           Df Deviance Resid. Df Resid. Dev  Pr(>Chi)
12 ## NULL                        59      82.911
13 ## supp           1      4.339      58      78.572    0.03724 *
14 ## factor(dose)   2     46.658      56      31.913 7.384e-11 ***
15 ## ---
16 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

1  exp(coef(fit)) # Exponentiate coefficients
2  ##      (Intercept)      suppVC factor(dose)1 factor(dose)2
3  ##  1.057677e-01  4.129916e-02  2.300199e+01  2.164477e+03
4  exp(confint(fit)) # 95% CI for Exponentiated coefficients
5  ##                2.5 %      97.5 %
6  ## (Intercept)  5.751916e-03 5.559953e-01
7  ## suppVC      1.843523e-03 3.157606e-01
8  ## factor(dose)1 2.832232e+00 5.260241e+02
9  ## factor(dose)2 1.010501e+02 1.701531e+05
10 drop1(fit, ~.      , test="Chisq") # Compare reduction in deviance,
    marginally
11 ## Single term deletions
12 ##
13 ## Model:
14 ## y ~ supp + factor(dose)
15 ##           Df Deviance      AIC      LRT  Pr(>Chi)
16 ## <none>           31.913 39.913
17 ## supp           1   42.802 48.802 10.888 0.0009678 ***

```

```
18  ## factor(dose)  2    78.572 82.572 46.658 7.384e-11 ***
19  ## ---
20  ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
21  confint(fit) # Confidence interval for the parameters
22  ##              2.5 %      97.5 %
23  ## (Intercept)  -5.158222 -0.5869954
24  ## suppVC      -6.296077 -1.1527709
25  ## factor(dose)1  1.041065  6.2653470
26  ## factor(dose)2  4.615616 12.0444538
```

## ◆ 回归表格输出

```
1 fit.1 = lm(len ~ factor(dose), ToothGrowth)
2 fit.2 = lm(len ~ factor(dose) + supp, ToothGrowth)
3
4 library(stargazer)
5 stargazer(fit.1, fit.2, fit, title="Regression Results",
6           type = 'html', style = 'asr')
```

Regression Results len y OLS logistic (1) (2) (3) factor(dose)1 9.130\*\*\* 9.130\*\*\* 3.136\*  
factor(dose)2 15.495\*\*\* 15.495\*\*\* 7.680\*\*\* suppVC -3.700\*\*\* -3.187\*\* Constant 10.605\*\*\*  
12.455\*\*\* -2.247\* N 60 60 60 R2 0.703 0.762 Adjusted R2 0.692 0.750 Log Likelihood  
-15.957 Residual Std. Error 4.242 (df = 57) 3.828 (df = 56) F Statistic 67.416\*\*\* (df = 2; 57)  
59.880\*\*\* (df = 3; 56) AIC 39.913  $p < .05$ ;  $p < .01$ ;  $p < .001$

## ◆ stargazer

Table 7: Mining Sector Expansion and Local Energy Prices

	(1) Average All Use Gas	(2) Industrial Gas Use	(3) Electricity All Use
<i>Instrumental Variables:</i>			
Mining Sector Share	-2.409* (1.262)	-6.067** (2.453)	-0.188* (0.113)
<i>Reduced Form:</i>			
Shale x Post 2008	-0.022** (0.011)	-0.055*** (0.019)	-0.002* (0.001)
<i>Ordinary Least Squares:</i>			
Mining Sector Share	-0.059 (0.069)	-0.033 (0.149)	-0.003 (0.005)
Clusters	337	337	364
Observations	24187	24620	33849
Instrument	15.09	14.99	15.79
R-squared	.953	.898	.885

Notes: All regressions include state-time fixed effects and county fixed effects. Column (1) uses the log of average gas prices in a county, where the consumer, commercial and industrial gas prices are weighted by their national consumption shares. In column (2) I only study the price charged to industrial consumers. Column (3) is the level of average electricity prices, where consumer, commercial and industrial prices are weighted by their respective national consumption shares. Robust standard errors clustered at the workforce investment board area are given in the parentheses with stars indicating \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .