

# STAT R note

Wenxiao Yang\*

\*Department of Mathematics, University of Illinois at Urbana-Champaign

2021

## 目录

<b>1 Basic</b>	<b>5</b>
1.1 q-value of $\chi_n^2$	5
1.2 读取数据 txt (galton)	5
1.3 读取数据 csv (bikeshares)	5
1.4 查看数据维度 (bikeshares)	5
1.5 数据中删除列 (bikeshares)	5
1.6 数据“列”处理: 赋值, 条件选中 (galton)	5
1.7 查看数据类型 (bikeshare)	6
1.8 data.frame	6
1.8.1 修改列名	6
1.8.2 data.frame 列拼接 cbind() (bikeshares)	6
1.8.3 data.frame 行拼接 rbind() (bikeshares)	7
1.8.4 data.frame 抽样	7
1.9 集体求均值	7
1.10 numeric	8
1.10.1 numeric(k): 生成 k 个 0 的 numeric	8
1.10.2 numeric 数值修改	8
1.11 matrix	8
1.11.1 data.frame 转成 matrix	8
1.11.2 修改列名	8
1.11.3 去掉矩阵列/行的名字	8
1.11.4 自己创建 matrix	8
1.11.5 Transpose of matrix 转置矩阵	9
1.11.6 Multiplication of matrix 矩阵乘法	9
1.11.7 解 $Ax = b$ : solve(A,b)	9

1.11.8	矩阵行列式: <code>det()</code>	9
1.11.9	生成对角阵: <code>diag(1,2,3,4)</code>	10
1.11.10	提取对角线上的元素: <code>diag()</code>	10
1.11.11	特征值和特征向量: <code>eigen()</code>	10
1.11.12	逆矩阵 <code>solve(A)</code>	10
1.11.13	列或行的函数处理 <code>apply(A, 1/2, func)</code>	11
<b>2</b>	<b>Simple Linear Regression</b>	<b>11</b>
2.1	拟合 <code>slr (galton)</code>	11
2.2	Summary 中提取 R-square ( <code>galton</code> )	12
2.3	Summary 中提取 coefficients ( <code>galton</code> )	12
2.4	回归中提取 degrees of freedom ( <code>galton</code> )	12
2.5	Hypothesis test	12
2.5.1	p-value of t-test ( <code>galton</code> )	12
2.5.2	Critical value of $\alpha = 0.05$ in <code>t(n)</code>	12
2.5.3	ANOVA(F-test) ( <code>HW1</code> )	12
2.5.4	p-value of F-test ( <code>HW1</code> )	13
2.5.5	Critical value of $\alpha = 0.05$ in <code>F(p,n)</code>	13
2.6	Confidence interval 置信区间 ( <code>HW1</code> )	13
2.7	Prediction	13
2.7.1	模型带入数据 ( <code>galton</code> )	13
2.7.2	Confidence interval ( <code>HW1</code> ) $\hat{\beta}_0 + \hat{\beta}_1 x^* \pm T_{n-2}(\alpha/2) \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$	14
2.7.3	Prediction interval ( <code>HW1</code> ) $\hat{\beta}_0 + \hat{\beta}_1 x^* \pm T_{n-2}(\alpha/2) \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$	14
<b>3</b>	<b>Multiple Linear Regression</b>	<b>14</b>
3.1	拟合 <code>mlr (bikeshares)</code>	14
3.2	Update regression, add or delete predictor	15
3.3	回归中提取 residuals, fitted values ( <code>bikeshare</code> )	15
3.4	Summary 中提取 F-test statistic	15
3.4.1	得到 RSS: $\sum_{i=1}^n r_i^2$	15
3.5	Correlation matrix ( <code>bikeshares</code> ) <code>cor()</code>	15
3.6	Plot all pairs of variables	16
3.7	Partial F-Tests ( <code>bikeshare</code> )	16
3.8	Permutation Tests ( <code>bikeshares</code> )	17
3.9	Confidence/Prediction Interval	18
3.9.1	Estimators' Confidence Interval	18
3.9.2	Estimators' Confidence regions	18
3.9.3	Confidence Interval for new observation	19
3.9.4	Prediction Interval for new observation	20

3.10	Unusual Observation . . . . .	20
3.10.1	Leverage Points . . . . .	20
3.10.2	Half-norm Plot . . . . .	20
3.10.3	Standardized Residuals, Studentized Residuals, <i>rstandard()</i> , <i>rstudent()</i> . . . .	21
3.10.4	Outliers . . . . .	21
3.11	High influential points . . . . .	22
3.12	Diagnostics . . . . .	23
3.12.1	Checking Homoskedasticity Graph . . . . .	23
3.12.2	Breusch-Pagan Test . . . . .	24
3.12.3	Checking Normality Graph . . . . .	24
3.12.4	Shapiro test . . . . .	25
3.12.5	Kolmogorov-Smirnov test . . . . .	25
3.12.6	Checking Serial Dependence: Durbin Watson test . . . . .	25
3.12.7	Checking the Linearity Assumption with Partial Regression Plots . . . . .	26
3.12.8	Box Cox Transformations . . . . .	26
3.12.9	Summary of Diagnostic Plots . . . . .	28
3.13	Collinearity . . . . .	29
3.13.1	Standardized each colum of $X$ . . . . .	29
3.13.2	Condition number of the $\mathbf{X}^T\mathbf{X}$ matrix . . . . .	30
3.13.3	Variance Inflation Factor (VIF) . . . . .	30
3.13.4	Pairwise correlations and partial F-tests . . . . .	30
<b>4</b>	<b>Time Series</b>	<b>31</b>
4.1	First Order Autoregressive Model . . . . .	31
<b>5</b>	<b>Polynomials Regression</b>	<b>32</b>
5.1	Orthogonal Polynomials . . . . .	32
5.2	B-Splines Basis . . . . .	32
5.3	Natural Cubic Splines . . . . .	33
<b>6</b>	<b>Categorical ANOVA</b>	<b>33</b>
6.1	Effect tests . . . . .	33
6.2	ANOVA Type III . . . . .	33
<b>7</b>	<b>Variation Selection</b>	<b>34</b>
7.1	Leap and Bounds method . . . . .	34
7.2	Searching algorithm based on AIC and BIC . . . . .	34
<b>8</b>	<b>Shrinkage Methods</b>	<b>35</b>
8.1	PCR, PCA . . . . .	35

<b>9</b>	<b>One/Two Way ANOVA</b>	<b>36</b>
9.1	Pairwise comparisons . . . . .	36
<b>10</b>	<b>Experimental Designs</b>	<b>37</b>
10.1	Paired t-test . . . . .	37
10.2	. . . . .	37
<b>11</b>	<b>画图</b>	<b>37</b>
11.1	$2 \times 2$ 的画布 . . . . .	37
11.2	plot 点图, 接上节 (bikeshares) . . . . .	37
11.3	ggplot . . . . .	37
11.3.1	Plot the regression line along with the connected “point-wise” confidence inter- vals (galton) . . . . .	38
11.3.2	给颜色取名, 竖直的线, 坐标 label . . . . .	38

## 1 Basic

### 1.1 q-value of $\chi_n^2$

```
1 qchisq(0.95, n)
```

### 1.2 读取数据 txt (galton)

```
1 galton <- read.table("Galton.txt", header=TRUE)
```

### 1.3 读取数据 csv (bikeshares)

```
1 bikeshares <- read.csv("BikeShares.csv", header=TRUE)
2 #分隔符为";"
3 whitewines.data<-read.csv("whitewines.csv", sep=";", header = TRUE)
```

### 1.4 查看数据维度 (bikeshares)

```
1 dim(bikeshares)
2 ## [1] 17414 10
```

### 1.5 数据中删除列 (bikeshares)

```
1 # We remove columns 1, 7, 8, 9, 10:
2 bikeshares.reg = bikeshares[,c(-1,-7,-8,-9,-10)] # -i 即删除 i 列
3 head(bikeshares.reg)
4 ## cnt t1 t2 hum wind_speed
5 ## 1 182 3.0 2.0 93.0 6.0
6 ## 2 138 3.0 2.5 93.0 5.0
7 ## 3 134 2.5 2.5 96.5 0.0
8 ## 4 72 2.0 2.0 100.0 0.0
9 ## 5 47 2.0 0.0 93.0 6.5
10 ## 6 46 2.0 2.0 93.0 4.0
```

### 1.6 数据“列”处理：赋值，条件选中 (galton)

```

1 # Define the Adjusted Height Variable (according to Galton)
2 galton$AH <- galton$Height
3 galton$AH[galton$Gender=="F"]<-galton$Height[galton$Gender=="F"]*1.08
4 galton$MP <- (galton$Father + 1.08*galton$Mother)/2
5 head(galton)
6 ##   Family Father Mother Gender Height Kids    AH    MP
7 ## 1      1   78.5   67.0      M   73.2    4 73.200 75.43
8 ## 2      1   78.5   67.0      F   69.2    4 74.736 75.43
9 ## 3      1   78.5   67.0      F   69.0    4 74.520 75.43
10 ## 4      1   78.5   67.0      F   69.0    4 74.520 75.43
11 ## 5      2   75.5   66.5      M   73.5    4 73.500 73.66
12 ## 6      2   75.5   66.5      M   72.5    4 72.500 73.66

```

## 1.7 查看数据类型 (bikeshare)

```

1 class(numeric(n.iter))
2 ## [1] "numeric"
3 class(bikeshares.reg)
4 ## [1] "data.frame"

```

## 1.8 data.frame

### 1.8.1 修改列名

```

1 names(myCR) = c("t1", "hum");

```

### 1.8.2 data.frame 列拼接 cbind() (bikeshares)

```

1 bikeshare.mlr1$fitted[1:5]
2 ## 1      2      3      4      5
3 ## 158.12967 152.85747 42.50091 -77.95731 126.47427
4 bikeshare.mlr1$residuals[1:5]
5 ## 1      2      3      4      5
6 ## 23.87033 -14.85747 91.49909 149.95731 -79.47427
7 cbind(bikeshare.mlr1$fitted[1:5], bikeshare.mlr1$residuals[1:5])
8 ##      [,1]      [,2]
9 ## 1 158.12967 23.87033

```

```

10 ## 2 152.85747 -14.85747
11 ## 3 42.50091 91.49909
12 ## 4 -77.95731 149.95731
13 ## 5 126.47427 -79.47427

```

### 1.8.3 data.frame 行拼接 rbind() (bikeshares)

```

1 rbind(bikeshare.mlr1$fitted[1:5], bikeshare.mlr1$residuals[1:5])
2 ## 1 2 3 4 5
3 ## [1,] 158.12967 152.85747 42.50091 -77.95731 126.47427
4 ## [2,] 23.87033 -14.85747 91.49909 149.95731 -79.47427

```

### 1.8.4 data.frame 抽样

```

1 head(bikeshares.reg)
2 ## cnt t1 t2 hum wind_speed
3 ## 1 182 3.0 2.0 93.0 6.0
4 ## 2 138 3.0 2.5 93.0 5.0
5 ## 3 134 2.5 2.5 96.5 0.0
6 ## 4 72 2.0 2.0 100.0 0.0
7 ## 5 47 2.0 0.0 93.0 6.5
8 ## 6 46 2.0 2.0 93.0 4.0
9 bikeshares.reg[sample(5), c(3,4)] #前五（第3，4列）中随机抽样
10 ## t2 hum
11 ## 2 2.5 93.0
12 ## 5 0.0 93.0
13 ## 3 2.5 96.5
14 ## 1 2.0 93.0
15 ## 4 2.0 100.0

```

## 1.9 集体求均值

```

1 apply(bikeshares.reg, 2, mean)
2 ## cnt t1 t2 hum wind_speed
3 ## 1143.10164 12.46809 11.52084 72.32495 15.91306

```

## 1.10 numeric

### 1.10.1 numeric(k): 生成 k 个 0 的 numeric

```
1 numeric(5)
2 ## [1] 0 0 0 0 0
3 class(numeric(5))
4 ## [1] "numeric"
```

### 1.10.2 numeric 数值修改

```
1 A=numeric(5)
2 A[1]=2
3 A
4 ## [1] 2 0 0 0 0
```

## 1.11 matrix

### 1.11.1 data.frame 转成 matrix

```
1 M=data.matrix(X)
```

### 1.11.2 修改列名

```
1 colnames(x)=c("t1", "t2", "hum")
```

### 1.11.3 去掉矩阵列/行的名字

```
1 rownames(A)<-NULL
2 colnames(A)<-NULL
```

### 1.11.4 自己创建 matrix

```
1 A=matrix(1:12,nrow=3,ncol=4)
2 A
3 ##      [,1] [,2] [,3] [,4]
4 ## [1,]    1    4    7   10
```



```

5 ## [2 ,]      2      5      8     11
6 ## [3 ,]      3      6      9     12

```

### 1.11.5 Transpose of matrix 转置矩阵

```

1 t(A)
2 ##           [,1] [,2] [,3]
3 ## [1 ,]      1      2      3
4 ## [2 ,]      4      5      6
5 ## [3 ,]      7      8      9
6 ## [4 ,]     10     11     12

```

### 1.11.6 Multiplication of matrix 矩阵乘法

```

1 A%*%t(A)
2 ##           [,1] [,2] [,3]
3 ## [1 ,]     166    188    210
4 ## [2 ,]     188    214    240
5 ## [3 ,]     210    240    270

```

### 1.11.7 解 $Ax = b$ : solve(A,b)

Solve  $ax = b$

```

1 A.data=data.frame(a=c(1,43,765,9),b=c(2,455,787,2),
2                   c=c(213,434,67,24),d=c(672,332,7,123))
3 A=data.matrix(A.data)
4 b=matrix(c(1,10,8,9))
5 solve(A,b)
6 ##           [,1]
7 ## a  -0.6723499
8 ## b   0.7380811
9 ## c  -0.9034118
10 ## d   0.2866412

```

### 1.11.8 矩阵行列式: det()

```
1 det(A)
```

### 1.11.9 生成对角阵: `diag(1,2,3,4)`

```
1 diag(c(1,2,3,4))
2 ##      [,1] [,2] [,3] [,4]
3 ## [1,]    1    0    0    0
4 ## [2,]    0    2    0    0
5 ## [3,]    0    0    3    0
6 ## [4,]    0    0    0    4
```

### 1.11.10 提取对角线上的元素: `diag()`

```
1 diag(A)
2 ## [1]    1 455    67 123
```

### 1.11.11 特征值和特征向量: `eigen()`

```
1 eigen(A)
2 ## eigen() decomposition
3 ## $values
4 ## [1]  962.54862 -533.15335  195.96895   20.63578
5 ##
6 ## $vectors
7 ##      [,1]      [,2]      [,3]      [,4]
8 ## [1,] -0.18050353 -0.31476395  0.7098847  0.5218457
9 ## [2,] -0.65689212 -0.36245740 -0.6561850 -0.5428550
10 ## [3,] -0.73165231  0.87683413  0.2141936  0.6319310
11 ## [4,] -0.02441547 -0.02664961  0.1400217 -0.1834356
```

### 1.11.12 逆矩阵 `solve(A)`

```
1 solve(A)
2 ##      [,1]      [,2]      [,3]      [,4]
3 ## [1,]  0.015470466 -0.0038533021  0.0023771584 -0.07425607
4 ## [2,] -0.016656510  0.0038972675 -0.0011449021  0.08054712
```

```

5 ## [3,]  0.019498924 -0.0018420827  0.0012737127 -0.10163107
6 ## [4,] -0.004665816  0.0005780095 -0.0004038514  0.03208420

```

### 1.11.13 列或行的函数处理 $apply(A, 1/2, func)$

```

1 apply(A,1,mean) #1 表示对行求均值
2 apply(A,2,mean) #2 表示对列求均值
3 apply(x,2,sd)
4 apply(x,2,var)

```

## 2 Simple Linear Regression

### 2.1 拟合 slr (galton)

```

1 # Simple Linear Regression
2 slr.fit <- lm(AH ~ MP, data=galton)
3 summary(slr.fit)
4 ##
5 ## Call:
6 ## lm(formula = AH ~ MP, data = galton)
7 ##
8 ## Residuals:
9 ##      Min       1Q   Median       3Q      Max
10 ## -9.4947 -1.4779  0.0995  1.5175  9.1262
11 ##
12 ## Coefficients:
13 ##              Estimate Std. Error t value Pr(>|t|)
14 ## (Intercept) 18.76698     2.84062   6.607 6.74e-11 ***
15 ## MP          0.72906     0.04102  17.772 < 2e-16 ***
16 ## ---
17 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
18 ##
19 ## Residual standard error: 2.233 on 896 degrees of freedom
20 ## Multiple R-squared:  0.2606, Adjusted R-squared:  0.2598
21 ## F-statistic: 315.9 on 1 and 896 DF, p-value: < 2.2e-16

```

## 2.2 Summary 中提取 R-square (galton)

```
1 summary(slr.fit)$r.square
```

## 2.3 Summary 中提取 coefficients (galton)

```
1 galton.coef = summary(slr.fit)$coef
2 galton.coef
3 ##              Estimate Std. Error   t value    Pr(>|t|)
4 ## (Intercept) 18.7669821  2.84062068   6.606648 6.735528e-11
5 ## MP          0.7290562  0.04102226  17.772211 9.224505e-61
6
7 galton.coef[2,1]
8 galton.coef[2,3] ## 提取 t-test
```

## 2.4 回归中提取 degrees of freedom (galton)

```
1 slr.fit$df
2 ## [1] 896
```

## 2.5 Hypothesis test

### 2.5.1 p-value of t-test (galton)

```
1 # pt(t-statistics, df)
2 # $H_0: \beta_1=0$, 由于检验 0 对称, 我们需要乘 2
3 2*pt(-galton.coef[2,1]/galton.coef[2,2], 896)
4 ## [1] 9.224505e-61
```

### 2.5.2 Critical value of $\alpha = 0.05$ in $t(n)$

```
1 qt(.05, n)
2 ## -1.644941
```

### 2.5.3 ANOVA(F-test) (HW1)

```

1 grade.anova=anova(slr.fit)
2 grade.anova
3 ## Analysis of Variance Table
4 ##
5 ## Response: final
6 ##           Df SumSq MeanSq F value Pr(>F)
7 ## QuizAverage  1  69812   69812  423.19 < 2.2e-16 ***
8 ## Residuals   380  62687    165
9 ## ———
10 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
11
12 grade.anova[1,4] ## 提取F- value from ANOVA Table

```

#### 2.5.4 p-value of F-test (HW1)

```

1 pf(grade.anova[1,4], df1=1, df2=380, lower.tail = FALSE)
2 # lower.tail: if TRUE (default), probabilities are P[X ≤ x],
3 # otherwise, P[X > x].

```

#### 2.5.5 Critical value of $\alpha = 0.05$ in F(p,n)

```

1 qf(.05, p, n, lower.tail = FALSE)

```

#### 2.6 Confidence interval 置信区间 (HW1)

```

1 confint(slr.fit, 'QuizAverage', level=0.9)
2 ##           5 %           95 %
3 ##QuizAverage 0.7880018 0.9253306

```

#### 2.7 Prediction

##### 2.7.1 模型带入数据 (galton)

```

1 predict(slr.fit, newdata=data.frame(MP=70))
2 ##           1
3 ## 69.80092

```

### 2.7.2 Confidence interval (HW1) $\hat{\beta}_0 + \hat{\beta}_1 x^* \pm T_{n-2}(\alpha/2) \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$

```

1 predict(slr.fit, newdata = data.frame(QuizAverage=85),
2 interval = 'confidence', level=0.9)
3 ##          fit          lwr          upr
4 ## 1 76.7638 75.44682 78.08077

```

### 2.7.3 Prediction interval (HW1) $\hat{\beta}_0 + \hat{\beta}_1 x^* \pm T_{n-2}(\alpha/2) \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$

```

1 predict(slr.fit, newdata = data.frame(QuizAverage=85),
2 interval = 'prediction', level=0.9)
3 ##          fit          lwr          upr
4 ## 1 76.7638 55.54486 97.98273

```

## 3 Multiple Linear Regression

### 3.1 拟合 mlr (bikeshares)

```

1 bikeshare.mlr1 = lm(cnt ~ t1 + t2 + hum + wind_speed,
2                      data=bikeshares.reg)
3 summary(bikeshare.mlr1)
4 ##
5 ## Call:
6 ## lm(formula = cnt ~ t1 + t2 + hum + wind_speed,
7 ##      data = bikeshares.reg)
8 ##
9 ## Residuals:
10 ##      Min       1Q   Median       3Q      Max
11 ## -1970.1  -602.7  -252.7   332.6  6007.4
12 ##
13 ## Coefficients:
14 ##              Estimate Std. Error t value Pr(>|t|)
15 ## (Intercept) 2582.5618    64.7237   39.901 < 2e-16 ***
16 ## t1           66.1963     9.4206    7.027 2.19e-12 ***
17 ## t2          -18.2313     7.7565   -2.350 0.018762 *
18 ## hum         -27.5645     0.5865 -46.999 < 2e-16 ***
19 ## wind_speed   -3.8435     0.9899   -3.883 0.000104 ***

```

```

20 ## —
21 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
22 ##
23 ## Residual standard error: 936 on 17409 degrees of freedom
24 ## Multiple R-squared:  0.2562, Adjusted R-squared:  0.256
25 ## F-statistic: 1499 on 4 and 17409 DF,  p-value: < 2.2e-16

```

### 3.2 Update regression, add or delete predictor

```

1 rat.lm_body = update(rat.lm, ~ liver+dose)
2 rat.lm_body = lm(Y~liver+dose, data = rat)
3 # 两者等价

```

### 3.3 回归中提取 residuals, fitted values (bikeshare)

```

1 bikeshare.mlr1$res
2 bikeshare.mlr$residuals
3 bikeshare.mlr$fitted.values

```

### 3.4 Summary 中提取 F-test statistic

```

1 summary(bikeshare.mlr1)$fstat
2 ## value      numdf      dendif
3 ## 1499.07    4.00      17409.00
4 summary(bikeshare.mlr1)$fstat[1]
5 ## 1499.07

```

#### 3.4.1 得到 RSS: $\sum_{i=1}^n r_i^2$

```

1 sum(bikeshare.mlr1$res^2) #方法1
2 deviance(bikeshare.mlr1) #方法2

```

### 3.5 Correlation matrix (bikeshares) cor()

```

1 cor(bikeshares.reg[, -1]) #这里[, -1] 是不想算第一列
2 ##                      t1          t2          hum  wind_speed

```

```

3 ## t1          1.0000000  0.98834422 -0.4477810  0.14547097
4 ## t2          0.9883442  1.00000000 -0.4034951  0.08840854
5 ## hum         -0.4477810 -0.40349514  1.0000000 -0.28778917
6 ## wind_speed  0.1454710  0.08840854 -0.2877892  1.00000000

```

```

1 round(cor(seatpos), dig=2)
2 # 打印出来的数据保留两位小数

```

### 3.6 Plot all pairs of variables

```

1 pairs(rat)

```

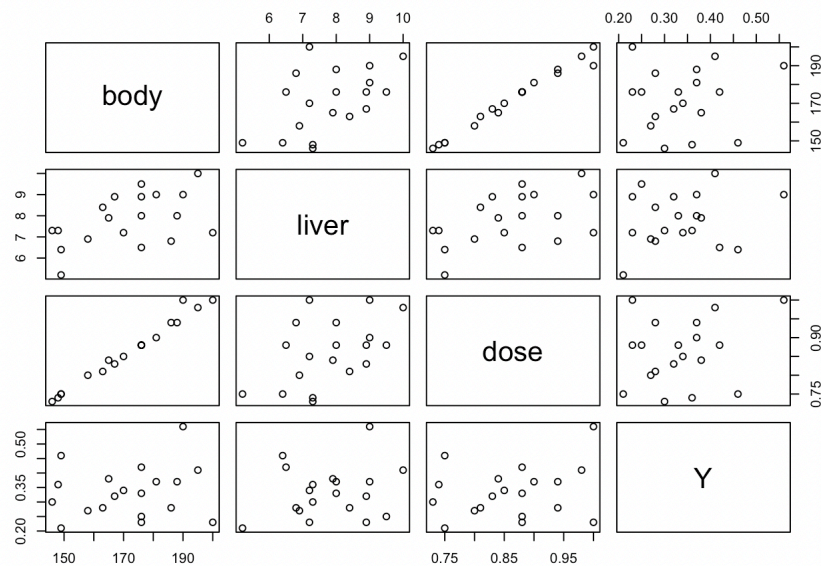


图 1:

### 3.7 Partial $F$ -Tests (bikeshare)

```

1 bikeshare.mlr.full = lm(cnt ~ t1 + t2+ hum + wind_speed ,
2                           data=bikeshares.reg ) #先回归 full model
3 bikeshare.mlr.reduced = lm(cnt ~ hum + wind_speed ,
4                             data=bikeshares.reg ) #回归 reduced model
5 anova(bikeshare.mlr.reduced , bikeshare.mlr.full)
6                                     #do the partial F-test by "anova(.)"

```



```

7  ## Analysis of Variance Table
8  ##
9  ## Model 1: cnt ~ hum + wind_speed
10 ## Model 2: cnt ~ t1 + t2 + hum + wind_speed
11 ##   Res.Df      RSS Df Sum of Sq      F    Pr(>F)
12 ## 1   17411 1.6103e+10
13 ## 2   17409 1.5250e+10  2 853010396 486.88 < 2.2e-16 ***
14 ## ———
15 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Sum of Square 853010396 是  $RSS_0 - RSS_\alpha = 1.6103e + 10 - 1.5250e + 10 = 853010396$

我们也可以按照公式算：

```

1  rss.full = sum(bikeshare.mlr.full$res^2)
2  # You can also compute it with
3  # deviance(bikeshare.mlr.full)
4  rss.reduced = sum(bikeshare.mlr.reduced$res^2)
5  # deviance(bikeshare.mlr.reduced)
6  Fstat = (rss.reduced - rss.full)/2/(rss.full/17409)
7  Fstat
8  ## [1] 486.8763
9  1-pf(Fstat, 2, 17409)
10 ## [1] 0

```

### 3.8 Permutation Tests (bikeshares)

$$\begin{cases} H_0 : \text{bikeshares} \sim \text{humidity} + \text{windspeed} \\ H_\alpha : \text{bikeshares} \sim \text{RealTemp} + \text{FeelsLikeTemp} + \text{humidity} + \text{windspeed} \end{cases}$$

If *RealTemp* and *FeelsLikeTemp* are insignificant (Under  $H_0$ ), the F-statistic of regression model will not be affected by switching the orders of these two data. Then new F-statistic will be equal(or less) to the old. i.e. High new F-statistic is more extreme than  $H_0$ . So lower  $p$ -value will support  $H_\alpha$  : *RealTemp* and *FeelsLikeTemp* are significant.

```

1  n.iter = 2000;
2  fstats = numeric(n.iter);
3  for(i in 1:n.iter){
4    newbikes = bikeshares.reg;
5    newbikes[, c(3,4)] = bikeshares.reg[sample(17414), c(3,4)];

```

```

6   ge = lm(cnt ~ t1 + t2+ hum + wind_speed, data=newbikes);
7   fstats[i] = summary(ge)$fstat[1]
8 }
9
10 # Estimated p-value
11 length(fstats[fstats > summary(bikeshare.mlr.full)$fstat[1]])/n.iter
12 ## [1] 0

```

## 3.9 Confidence/Prediction Interval

### 3.9.1 Estimators' Confidence Interval

```

1 confint(bikeshare.mlr)
2 ##                2.5 %      97.5 %
3 ## (Intercept) 2543.679114 2766.669772
4 ## t1          41.516598   47.099480
5 ## hum         -28.984942  -26.739780
6 ## wind_speed  -4.941603   -1.262794
7 confint(bikeshare.mlr, 't1', level=0.99)
8 ##      0.5 %   99.5 %
9 ## t1 40.63932 47.97676

```

### 3.9.2 Estimators' Confidence regions

```

1 library(ellipse)
2 library(ggplot2)
3 CR95 = ellipse(bikeshare.mlr, c(2,3))
4 CR99 = ellipse(bikeshare.mlr, c(2,3), level=0.99)
5 CR998 = ellipse(bikeshare.mlr, c(2,3), level=0.998)
6 # Plot Confidence Regions for column 2,3
7 dim(CR95)
8 ## [1] 100    2
9 head(CR95)
10 ##           t1      hum
11 ## [1,] 47.25426 -26.67754
12 ## [2,] 47.13012 -26.63239
13 ## [3,] 46.99462 -26.59219
14 ## [4,] 46.84830 -26.55710

```

```

15 ## [5,] 46.69175 -26.52728
16 ## [6,] 46.52561 -26.50282

```

```

1 myCR = rbind(CR95, CR99, CR998);
2 # 行连接
3 myCR = data.frame(myCR);
4 names(myCR) = c("t1", "hum");
5 myCR[, 'level'] = as.factor(c(rep(0.95, dim(CR95)[1]),
6                               rep(0.99, dim(CR99)[1]),
7                               rep(0.998, dim(CR998)[1])));
8 # 添加列 'level', 给各行根据精度赋值
9
10 ggplot(data=myCR, aes(x=t1, y=hum, colour=level)) +
11   geom_path(aes(linetype=level), size=1.5) +
12   geom_point(x=coef(bikeshare.mlr)[2], y=coef(bikeshare.mlr)[3]
13             , shape=3, size=3, colour='red') +
14   geom_point(x=0, y=0, shape=1, size=3, colour='red')

```

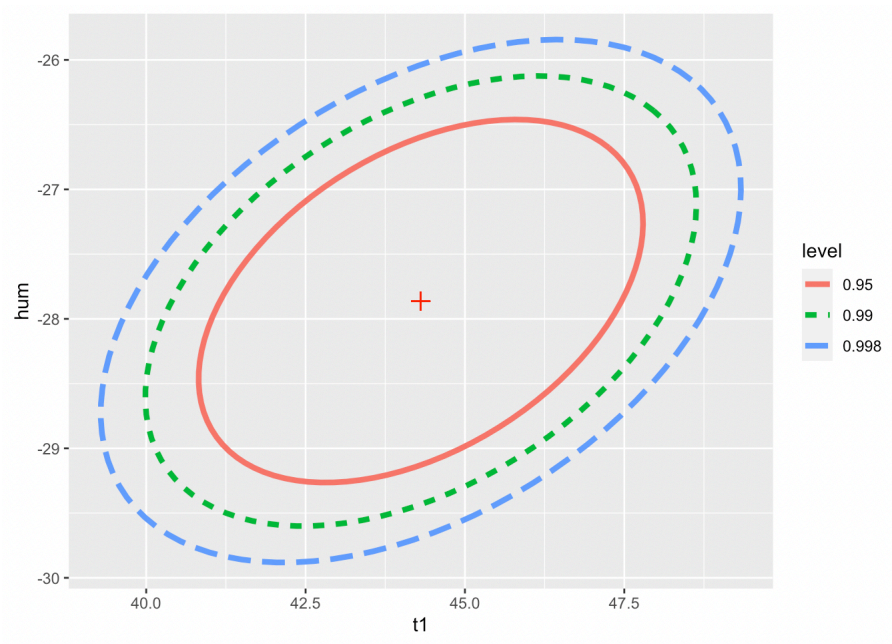


图 2:

### 3.9.3 Confidence Interval for new observation

```

1 x=data.frame(t(meanvalue))
2 predict.lm(bikeshare.mlr,x,interval="confidence",level=0.95)
3 ##          fit          lwr          upr
4 ## 1 1143.102 1129.198 1157.006

```

### 3.9.4 Prediction Interval for new observation

```

1 predict.lm(bikeshare.mlr,x,interval="prediction",level=0.95)
2 ##          fit          lwr          upr
3 ## 1 1143.102 -691.7461 2977.949

```

## 3.10 Unusual Observation

### 3.10.1 Leverage Points

```

1 lev=influence(bikeshare.mlr)$hat
2 # H matrix 的对角上的所有元素
3 newlev = lev[lev>2*p/n]
4 # 找出所有high leverage points
5 bikeshares.reg[lev > 2*p/n,]
6 # 筛选出bikeshares中high leverage points的项

```

### 3.10.2 Half-norm Plot

Designed to identify unusually large values and assess positive data.

Plot the data against the positive normal quantiles. Specifically,

1. Sort the data:

$$x_{[1]} \leq \dots \leq x_{[n]}.$$

2. Compute the quantiles:

$$u_i = \Phi^{-1}\left(\frac{n+i}{2n+1}\right)$$

3. Plot  $x_{[i]}$  against  $u_i$ .

```

1 library(faraway)
2 halfnorm(newlev, 6, labs=as.character(1:length(newlev)),
3   ylab="Leverages")
4 # 6是nlab, 即给几个点标注

```

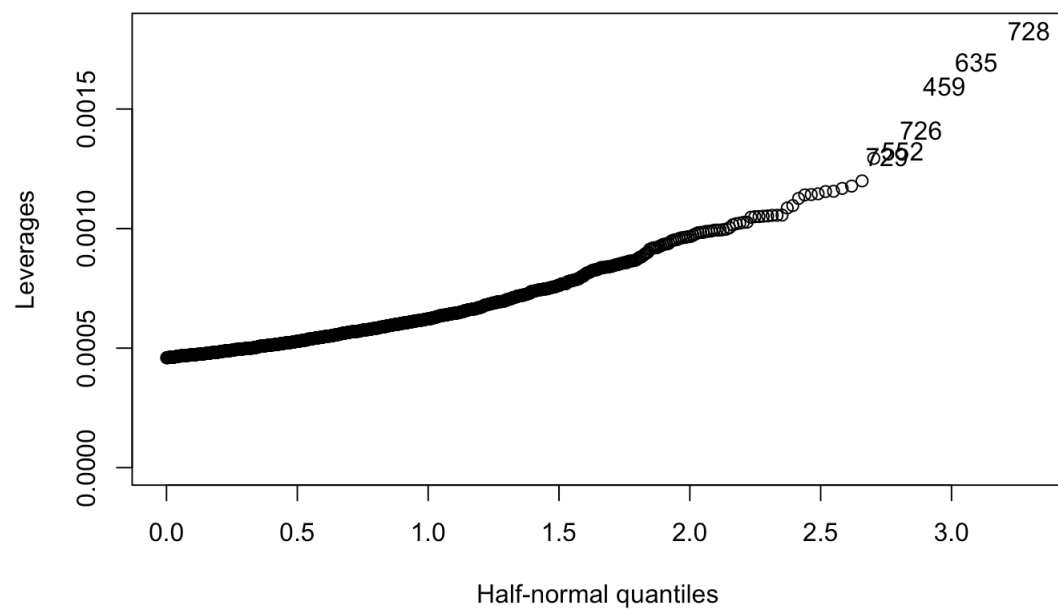


图 3:

### 3.10.3 Standardized Residuals, Studentized Residuals, *rstandard()*, *rstudent()*

```
1 rstandard(model)
2 rstudent(model)
```

### 3.10.4 Outliers

```
1 # Compute Studentized Residuals
2 jack=rstudent(bikeshare.mlr);
3 # The critical value WITH Bonferroni correction is
4 qt(.05/(2*n), n-p-1)
5 ## [1] -4.681361
6 # The critical value WITHOUT Bonferroni correction is
7 qt(.05/2, n-p-1)
8 ## [1] -1.9601
9 # Sort the residuals indescending order to find outliers (if any)
10 sort(abs(jack), decreasing=TRUE)[1:10]
11 ##      4462      5130      5139      4471      15888      5140
12 ## 6.408782 5.665958 5.499140 5.317999 4.807279 4.787554
13 ##      15217      15385      16727      14905
```

```
14 ## 4.746059 4.738005 4.661289 4.522918
```

As we can see here, we have 8 outliers, i.e. the values that are higher (in absolute value) of the critical T distribution value with Bonferroni correction ( $| - 4.681361 |$ ). These are observations: #4462, #5130, #5139, #4471, #15888, #5140, #15217, #15385.

### 3.11 High influential points

```
1 # Compute Cook's Distance
2 cook = cooks.distance(bikeshare.mlr)
3 # Extract max Cook's Distance
4 max(cook)
5 ## [1] 0.005641587
6 which.max(cook)
7 ## 4471
8 # Prepare a Half Normal Plot of Cook's distances
9 halfnorm(cook, 6, labs=as.character(1:length(cook)),
10 ylab="Cook's distances")
```

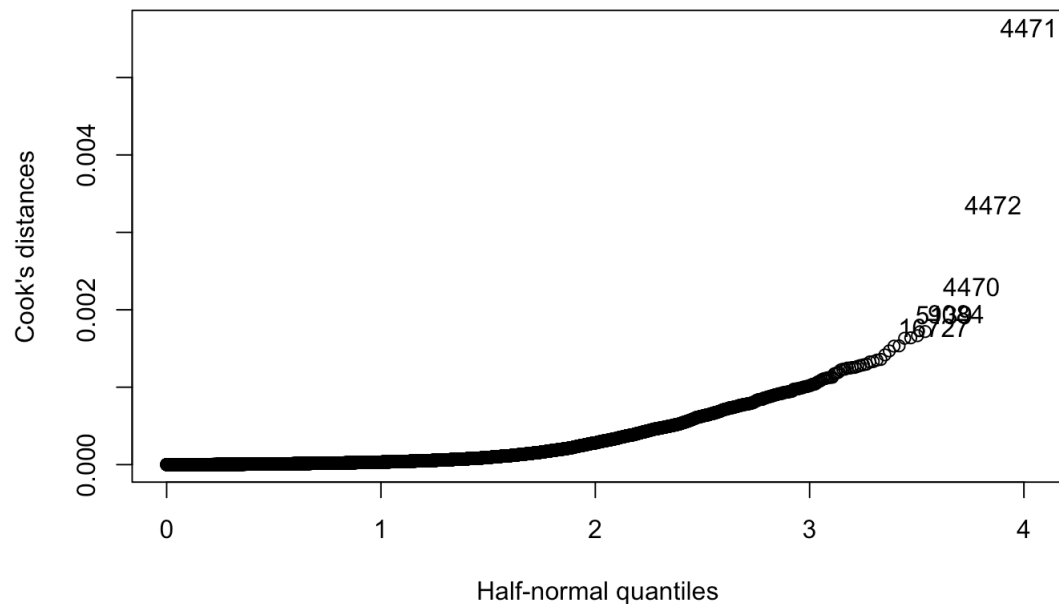


图 4:

## 3.12 Diagnostics

### 3.12.1 Checking Homoskedasticity Graph

```
1 plot(bikeshare.mlr, which=1)
```

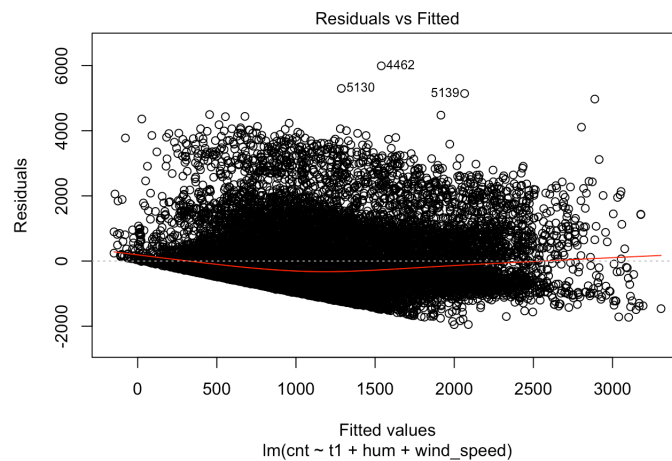


图 5:

Which is same as

```
1 plot(bikeshare.mlr$fitted.values, bikeshare.mlr$residuals)
```

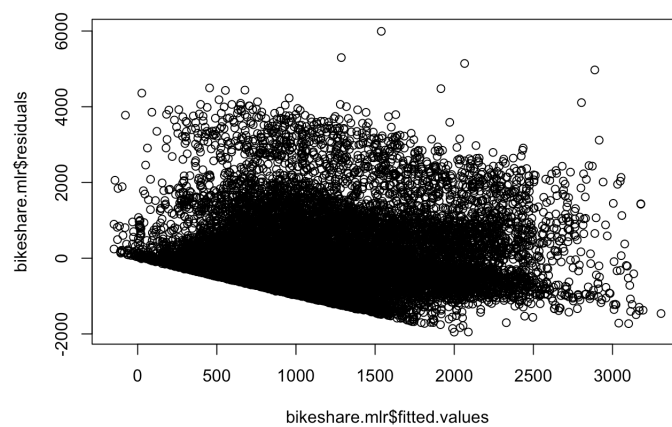


图 6:

### 3.12.2 Breusch-Pagan Test

```
1 library(lmtest)
2 bptest(bikeshare.mlr)
3 ##
4 ## studentized Breusch-Pagan test
5 ##
6 ## data: bikeshare.mlr
7 ## BP = 133.29, df = 3, p-value < 2.2e-16
```

We can also perform the BP test by hand:

```
1 tmp.fit = lm(bikeshare.mlr$res^2 ~ t1 + hum + wind_speed,
2   data=bikeshares.reg )
3 summary(tmp.fit)$r.sq*dim(bikeshares.reg)[1]
```

### 3.12.3 Checking Normality Graph

#### QQ-Plot

```
1 plot(bikeshare.mlr, which=2)
```

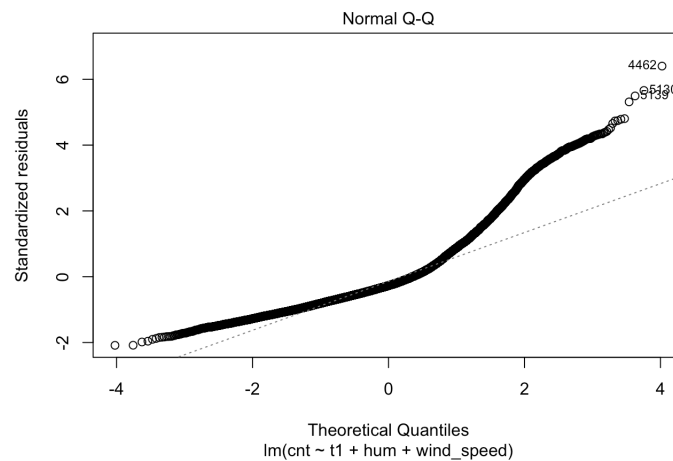


图 7:

#### Histogram



```
1 hist(bikeshare.mlr$residuals)
```

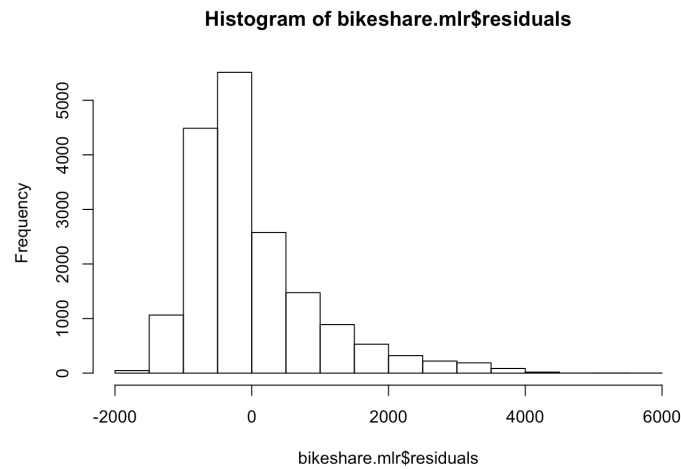


图 8:

### 3.12.4 Shapiro test

```
1 shapiro.test(residuals(bikeshare.mlr))
```

### 3.12.5 Kolmogorov-Smirnov test

```
1 ks.test(residuals(bikeshare.mlr), y=pnorm)
2 ##
3 ## One-sample Kolmogorov-Smirnov test
4 ##
5 ## data: residuals(bikeshare.mlr)
6 ## D = 0.63627, p-value < 2.2e-16
7 ## alternative hypothesis: two-sided
```

The p-value is low, which implies that the normality assumption is not satisfied either.

### 3.12.6 Checking Serial Dependence: Durbin Watson test

```
1 library(lmtest)
2 dwtest(lm.sales)
```

```

3 ##
4 ## Durbin-Watson test
5 ##
6 ## data: lm.sales
7 ## DW = 0.73473, p-value = 0.0001748
8 ## alternative hypothesis: true autocorrelation is greater than 0

```

### 3.12.7 Checking the Linearity Assumption with Partial Regression Plots

Test  $t_1$

```

1 bikeshare.mlr = lm(cnt ~ hum + wind_speed, data=bikeshares.reg)
2 bikeshare.mlr.t1 = lm(t1 ~ hum + wind_speed, data=bikeshares.reg)
3 plot(bikeshare.mlr.t1$residuals, bikeshare.mlr$residuals)

```

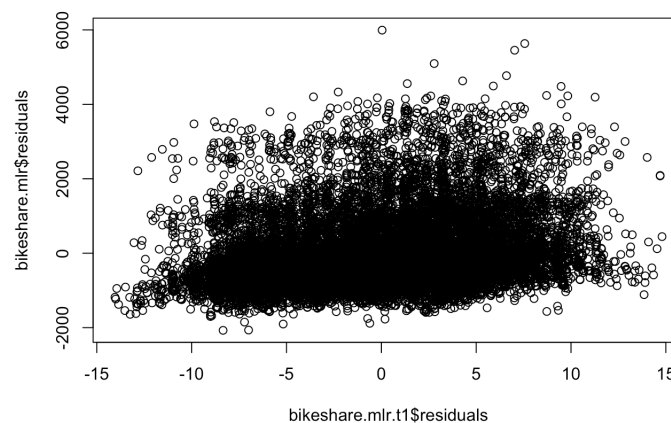


图 9:

### 3.12.8 Box Cox Transformations

First we need to make sure each  $y_i > 0$ :

```

1 min(bikeshares.reg$cnt) # this is the min value in the y's
2 ## [1] 0
3 which(bikeshares.reg$cnt==0)
4 # this is the location of the min value
5 ## [1] 2016

```

```

6 bikeshares.reg$cnt[2016]=0.01
7 # we replace the min with a small positive value
8 min(bikeshares.reg$cnt)
9 # we checke whether the 0 value was replaced
10 # by the small positive number
11 ## [1] 0.01

```

Now, we are ready to apply the *boxcox* function:

```

1 bikes.transformation = boxcox(bikeshare.mlr,
2 lambda=seq(-2, 2, length=400))

```

which also same as

```

1 boxcox(bikeshare.mlr, plotit=T) # plotit=T is the default setting

```

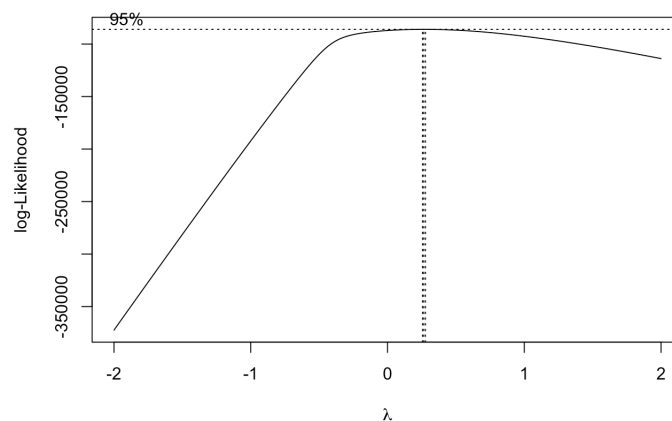


图 10:

改变范围

```

1 boxcox(bikeshare.mlr, plotit=T, lambda=seq(0.1, 1, by=0.05)) # zoom-in

```

Find the  $\lambda$  that maximizes the *Log – likelihood*.

```

1 names(bikes.transformation)

```

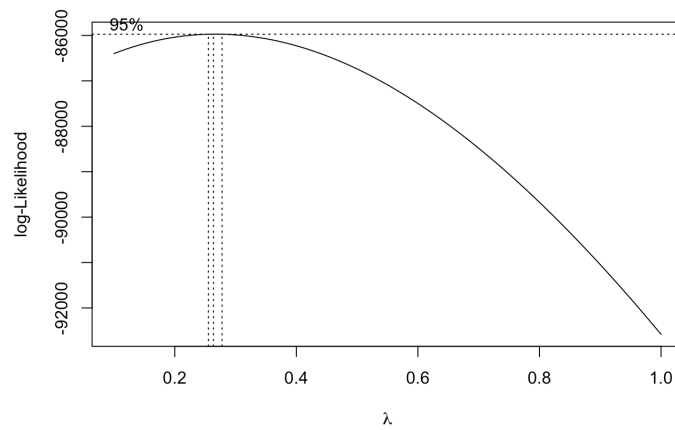


图 11:

```

2 ## [1] "x" "y"
3 bikes.transformation$x[1:10]
4 ## [1] -2.000000 -1.989975 -1.979950 -1.969925 -1.959900 -1.949875 -1.939850
5 ## [8] -1.929825 -1.919799 -1.909774
6 bikes.transformation$y[1:10]
7 ## [1] -372422.1 -370582.3 -368742.9 -366904.0 -365065.6 -363227.6 -361390.0
8 ## [8] -359553.0 -357716.4 -355880.2
9 bikes.transformation$x[bikes.transformation$y ==
10 max(bikes.transformation$y)] # lambda.hat
11 ## [1] 0.2656642

```

$$\hat{\lambda} = 0.2656642$$

Construct a Confidence Interval for  $\lambda$  as follows:

$$\{\lambda : L(\lambda) > L(\hat{\lambda}) - \frac{1}{2}\chi_1^2(1 - \alpha)\}$$

```

1 tmp=bikes.transformation$x[bikes.transformation$y >
2   max(bikes.transformation$y) - qchisq(0.95, 1)/2];
3 range(tmp) # 95% CI. Read Chapter 9 in the Faraway textbook for details.
4 ## [1] 0.2556391 0.2756892

```

### 3.12.9 Summary of Diagnostic Plots

```

1 fit=lm(Y~., data=rat)
2 par(mfrow=c(2,2))

```

```
3 plot(fit)
```

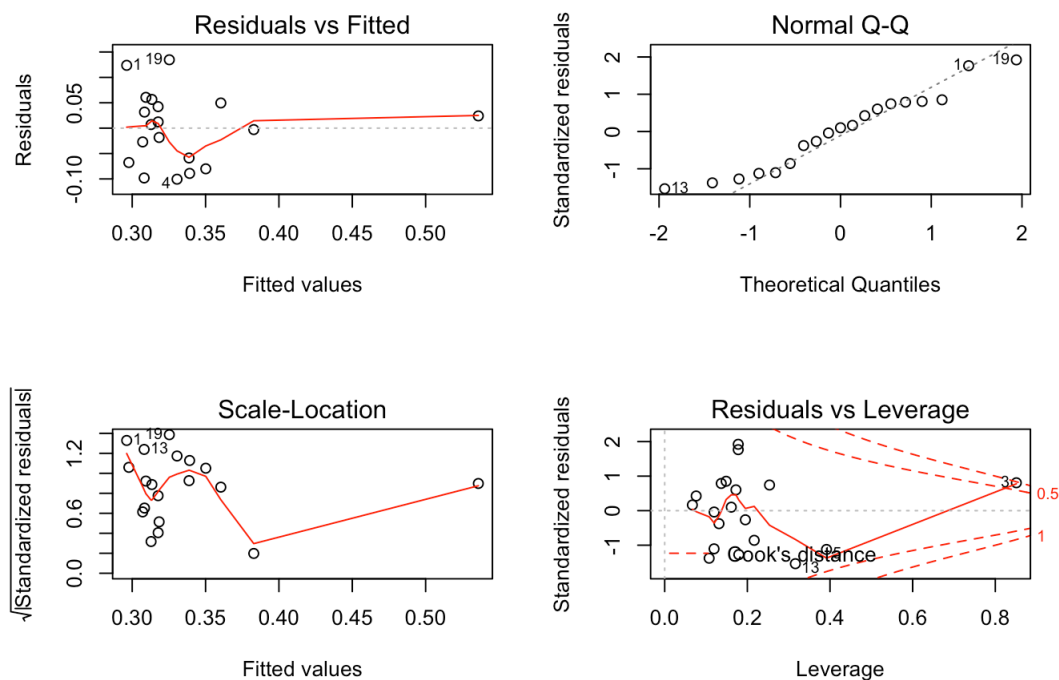


图 12:

### 3.13 Collinearity

```
1 library(faraway)
2 # 提取数据 seatpos
3 data(seatpos)
4 attach(seatpos)
5
6 # Fit the FULL model
7 position.full = lm(hipcenter ~ ., seatpos)
8 x = model.matrix(position.full)[, -1]
9 # remove the column that corresponds to the intercept
```

#### 3.13.1 Standardized each column of $X$

```

1 x = model.matrix(position.full)[,-1] #去除第一列 (即 intercept)
2 x = x - matrix(apply(x,2, mean), 38,8, byrow=TRUE)
3 x = x / matrix(apply(x, 2, sd), 38,8, byrow=TRUE)
4 apply(x,2,mean)
5 ##           Age           Weight           HtShoes           Ht           Seated
6 ## -2.193512e-17  2.810252e-16  9.566280e-16  1.941574e-16 -1.073010e-15
7 ##           Arm           Thigh           Leg
8 ## -1.070022e-16  8.909895e-17 -9.114182e-17
9 apply(x,2,var)
10 ##      Age  Weight HtShoes      Ht  Seated      Arm  Thigh      Leg
11 ##      1      1      1      1      1      1      1      1

```

### 3.13.2 Condition number of the $X^T X$ matrix

```

1 e = eigen(t(x) %*% x) # compute the eigenvalues
2 sqrt(e$val[1]/e$val[8])
3 ## [1] 59.7662

```

The condition number is 59.77, larger than 30, so we conclude that collinearity is present.

### 3.13.3 Variance Inflation Factor (VIF)

```

1 # Variance Inflation Factor (VIF)
2 round(vif(x), dig=2)
3 ##      Age  Weight HtShoes      Ht  Seated      Arm  Thigh      Leg
4 ##      2.00      3.65  307.43  333.14      8.95      4.50      2.76      6.69

```

```

1 sqrt(307.43)
2 ## [1] 17.53368

```

Note that the se for the coef associated with *HtShoes* is 17.5 times larger than it would have been without collinearity.

### 3.13.4 Pairwise correlations and partial F-tests

```

1 cor(Seated+Thigh, Ht)
2 ## [1] 0.9389819
3 cor(Seated+Leg, Ht)

```

```

4 ## [1] 0.965607
5 cor(Seated+Arm, Ht)
6 ## [1] 0.9465523

1 position.red1 = lm(hipcenter ~ Age + Weight + Ht + Seated, data=seatpos)
2 position.red2 = lm(hipcenter ~ Ht, data=seatpos)
3 anova(position.red2, position.red1)
4 ## Analysis of Variance Table
5 ##
6 ## Model 1: hipcenter ~ Ht
7 ## Model 2: hipcenter ~ Age + Weight + Ht + Seated
8 ##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
9 ## 1      36 47616
10 ## 2      33 44774   3    2841.6 0.6981 0.5599

```

Based on the  $F$ -test provided in the *ANOVA* table, we conclude that the reduced model with  $Ht$  as the only variable is better than the model that includes  $Age$ ,  $Weight$ ,  $Ht$  and  $Seated$ .

## 4 Time Series

### 4.1 First Order Autoregressive Model

```

1 library(nlme)
2 lm.sales.cor = gls(company_sales~industry_sales,
3 correlation = corAR1(form= ~ index), data=sales)
4
5 summary(lm.sales.cor)
6 ## Generalized least squares fit by REML
7 ##   Model: company_sales ~ industry_sales
8 ##   Data: sales
9 ##           AIC           BIC    logLik
10 ##   -31.74311  -28.18162  19.87156
11 ##
12 ## Correlation Structure: AR(1)
13 ##   Formula: ~index
14 ##   Parameter estimate(s):
15 ##   Phi
16 ##   1
17 ##

```

```

18 ## Coefficients:
19 ##
20 ## (Intercept)      -0.3189197  2041.6945  -0.00016   0.9999
21 ## industry_sales   0.1684878    0.0051  33.06272   0.0000
22 ##
23 ## Correlation:
24 ## (Intr)
25 ## industry_sales 0
26 ##
27 ## Standardized residuals:
28 ##      Min           Q1           Med           Q3
29 ## -9.036061e-05 -4.156415e-05 -3.013053e-06  8.080346e-05
30 ##  1.091922e-04
31 ##
32 ## Residual standard error: 2041.694
33 ## Degrees of freedom: 20 total; 18 residual

```

## 5 Polynomials Regression

### 5.1 Orthogonal Polynomials

```
1 poly(.)
```

### 5.2 B-Splines Basis

```

1 bs(x, df, knots, degree=3, intercept=FALSE)
2 # x是数据
3 # df是输出的 design matrix 的columns数, 和真正的df无关
4 # intercept=FALSE, df=真df-1
5 # intercept=TRUE, df=真df
6 # knots=k, 代表k是那个唯一的knot, 所以knot数是1, 无论k多大
7 new.knots= c(1/6, 3/6, 5/6)
8 bs(x, knots=new.knots, intercept=TRUE)
9 bs(x, knots=quantile(x, c(1/3,2/3)), intercept=TRUE)

```



## 5.3 Natural Cubic Splines

```
1 ns(x, df, knots, Boundary.knots, degree=3, intercept=FALSE)
2 # knots只表示interior knots, 还有俩boundary knots。
3 # 所以 真df==#knots+2
4 # 其他一样
5 ns(x, knots=new.knots, Boundary.knots=c(0,1), intercept=TRUE)
```

## 6 Categorical ANOVA

### 6.1 Effect tests

When the levels of the categorical variable are in text (instead of number), R assigns 0 and 1 in alphabetical order: 0 first and 1 second.

```
1 quest.full=lm(rate~lot.size*color, quest.data)
2 anova(quest.full)
3 ##Analysis of Variance Table
4 ##Response: rate
5 ##          Df Sum Sq Mean Sq F value    Pr(>F)
6 ##lot.size      1  43.226   43.226    7.1024  0.01765 *
7 ##color          1  20.052   20.052    3.2947  0.08955 .
8 ##lot.size:color  1   0.166    0.166    0.0272  0.87111
9 ##Residuals     15  91.293    6.086
10 ##——
11 ##Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

第一行 intercept only vs. intercept+lot.size

第二行 intercept+lot.size vs. intercept+lot.size+color

第三行 intercept+lot.size+color vs. intercept+lot.size+color+lot.size\*color

### 6.2 ANOVA Type III

This type tests for the presence of an effect given that both the other effects are in the model.

```
Anova(lm(1/time ~ treat*poison, data=newrats), type="III")
```

```
## Anova Table (Type III tests)
##
## Response: 1/time
##              Sum Sq Df F value    Pr(>F)
## (Intercept)  15.0605  1 66.5967 1.298e-09 ***
## treat         2.1340  3  3.1455  0.03723  *
## poison       11.7375  2 25.9514 1.225e-07 ***
## treat:poison  1.9800  6  1.4592  0.22073
## Residuals     7.9151 35
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

图 13:

## 7 Variation Selection

### 7.1 Leap and Bounds method

Use function *regsubsets* from library *leaps* to evaluate different scores for sub-sets of models up to size *p* (including the intercept).

```
1 library(leaps)
2 Hitters.leaps=regsubsets(Salary~., data=data.reg, nvmax=16)
3 rs=summary(Hitters.leaps)
4 rs$adjr2
5 rs$which[which.max(rs$adjr2),]
6 rs$cp
7 rs$bic
8 n=dim(data.reg)[1]
9 m=2:17
10 Aic=n*log(rs$rss/n)+2*m
```

### 7.2 Searching algorithm based on AIC and BIC

Use function *step* from the *stats* library to apply searching algorithms based on the AIC (default) or BIC criteria ( $k = \log(n)$ ). The option *direction=both* uses the Stepwise searching algorithm. You can also use the options: *direction = forward* and *direction = backward*.

```
1 step(full.model, direction="both")
2 step(full.model, direction="both", k=log(n))
3 We can also use direction=forward and direction=backward
```

## 8 Shrinkage Methods

### 8.1 PCR, PCA

Function `prcomp` can be used to calculate the PCs and extract the  $\lambda$ 's squared-roots (`sdev`) and eigenvectors (`rotation`) of the variance-covariance matrix:

```
data(meatspec, package="faraway")
trainmeat<-meatspec[1:172,]
testmeat<-meatspec[173:215,]
mod1<-lm(fat~., trainmeat)
meatpca<-prcomp(trainmeat[, -101])
round(meatpca$sdev, 3)[1:50]
```

```
## [1] 5.055 0.511 0.282 0.168 0.038 0.025 0.014 0.011 0.005 0.003 0.002 0.002
## [13] 0.001 0.001 0.001 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000
## [25] 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000
## [37] 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000
## [49] 0.000 0.000
```

图 14:

The `pcr` function (principal component regression) from the `pls` package has useful features for prediction and cross-validation. We can easily calculate the RMSE for the training set and the testing set.

```
1 library(pls)
2 modpcr<-pcr(fat ~ ., data=trainmeat, ncomp=50)
3 #summary(modpcr)
4 #RMSE with 4 PCAs
5 rmse(predict(modpcr, ncom=4), trainmeat$fat)
6 ## [1] 4.064745
7 rmse(predict(modpcr, testmeat, ncomp=4), testmeat$fat)
8 ## [1] 4.533982
```

You can use the function `RMSEP` instead, to select the number of PC's that minimize the 10-fold Cross-Validation error. The resulting Cross-Validation error is  $< 2.5$

```
1 set.seed(123)
2 # Minimize RMSE using function RMSEP
3 pcrmse<-RMSEP(modpcr, newdata=testmeat)
4 plot(pcrmse)
```

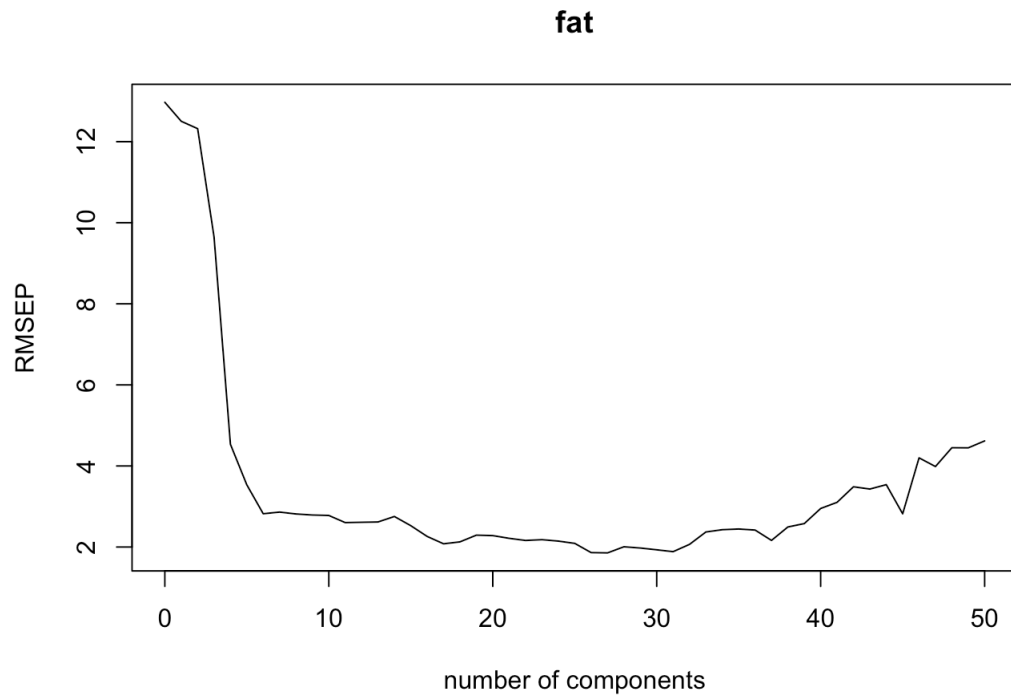


图 15:

## 9 One/Two Way ANOVA

### 9.1 Pairwise comparisons

Construct 90% family confidence intervals for all pairwise comparisons of classroom environments.

```
> TukeyHSD(aov(score ~ classroom, data=class.data), "classroom", conf.level=0.9)
Tukey multiple comparisons of means
90% family-wise confidence level

Fit: aov(formula = score ~ classroom, data = class.data)

$classroom
      diff      lwr      upr    p adj
B-A -3.0 -4.750606 -1.2493942 0.0029108
C-A -2.7 -4.450606 -0.9493942 0.0073313
C-B  0.3 -1.450606  2.0506058 0.9285599
```

图 16:

## 10 Experimental Designs

### 10.1 Paired t-test

```
1 t.test(shoes$A-shoes$B)
```

### 10.2

We use the `drop1` function instead of `anova`, because of the lack of orthogonality due to the incompleteness of the design.

```
1 lmodbibd <- lm(gain~block+treat, rabbit)
2 drop1(lmodbibd, test="F")
3 ## Single term deletions
4 ##
5 ## Model:
6 ## gain ~ block + treat
7 ##           Df Sum of Sq    RSS    AIC F value    Pr(>F)
8 ## <none>                150.77   78.437
9 ## block     9      595.74  746.51 108.426   6.5854 0.0007602 ***
10 ## treat     5      158.73  309.50  90.013   3.1583 0.0381655 *
11 ## ———
12 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

## 11 画图

### 11.1 $2 \times 2$ 的画布

```
1 par(mfrow=c(2,2))
```

### 11.2 plot 点图, 接上节 (bikeshares)

```
1 par(mfrow=c(2,2))
2 # Plot of t1 vs. cnt
3 plot(bikeshares.reg$t1, bikeshares.reg$cnt, xlab="Real
4 Temperature in C", ylab="New Bike Shares")
5 # Plot of t2 vs. cnt
6 plot(bikeshares.reg$t2, bikeshares.reg$cnt, xlab="Feels
```

```

7 Like_Temperature_in_C", ylab="New_Bike_Shares")
8 # Plot of t1 vs. t2
9 plot(bikeshares.reg$t1, bikeshares.reg$t2, xlab="Feels
10 Like_Temperature_in_C", ylab="Real_Temperature_in_C")
11 # Plot of hum vs. t1
12 plot(bikeshares.reg$hum, bikeshares.reg$t1, xlab="Humidity",
13 ylab="Real_Temperature_in_C")

```

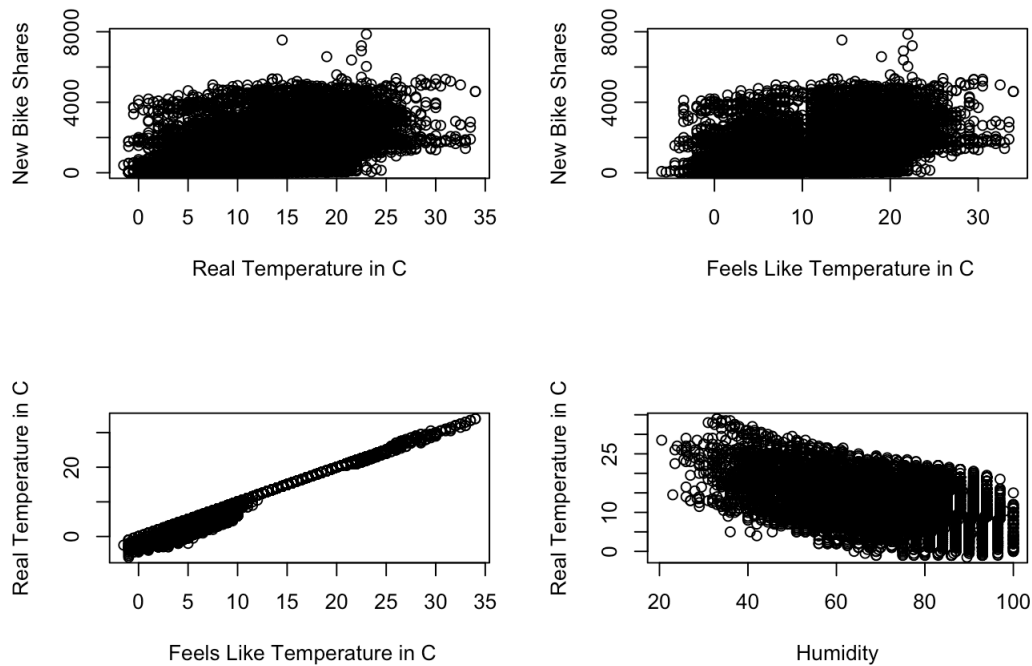


图 17:

### 11.3 ggplot

```

1 library(ggplot2)

```

#### 11.3.1 Plot the regression line along with the connected “point-wise” confidence intervals (galton)

```

1 library(ggplot2)

```

```
2 | ggplot(galton, aes(MP,AH)) + geom_point() + geom_smooth(method=lm)
```

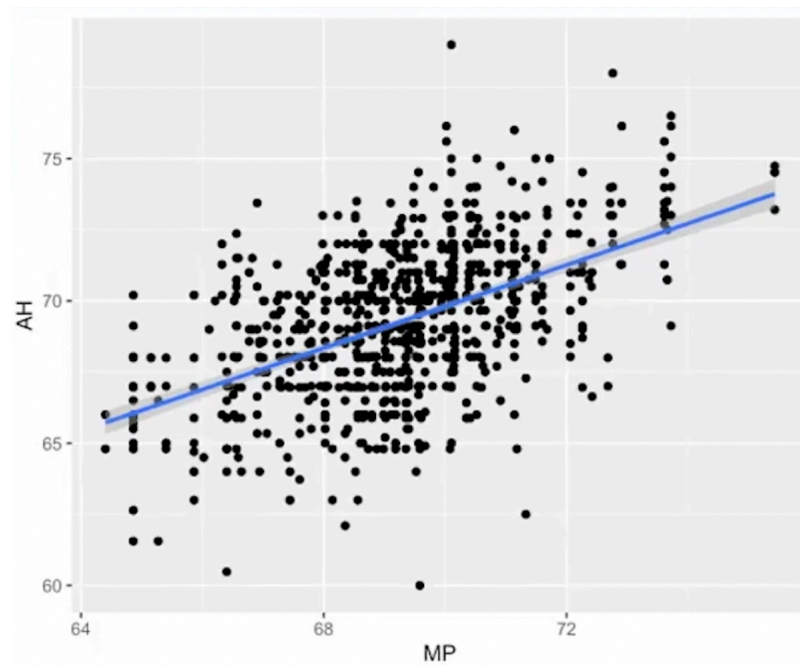


图 18:

### 11.3.2 给颜色取名，竖直的线，坐标 label

```
1 # Form the data frame for plotting
2 ggplot(data=NULL, aes(x=0:56)) +
3   geom_line(aes(y=myCI[,1], colour="LSfit"), size=1) +
4   geom_line(aes(y=myCI[,2], colour="90%_CI"), size=1) +
5   geom_line(aes(y=myCI[,3], colour="90%_CI"), size=1) +
6   geom_line(aes(y=myPI[,2], colour="90%_PI"), size=1, linetype=2)+
7   geom_line(aes(y=myPI[,3], colour="90%_PI"), size=1, linetype=2)+
8   scale_colour_manual("", values=c("LSfit" = "black",
9                                     "90%_CI" = "blue",
10                                    "90%_PI"="red"))+
11   xlab("wind_speed") +ylab("bike_shares")+
12   geom_vline(xintercept = mean(bikeshares.reg$wind_speed),
13             colour="purple", size=1, linetype=3)
```

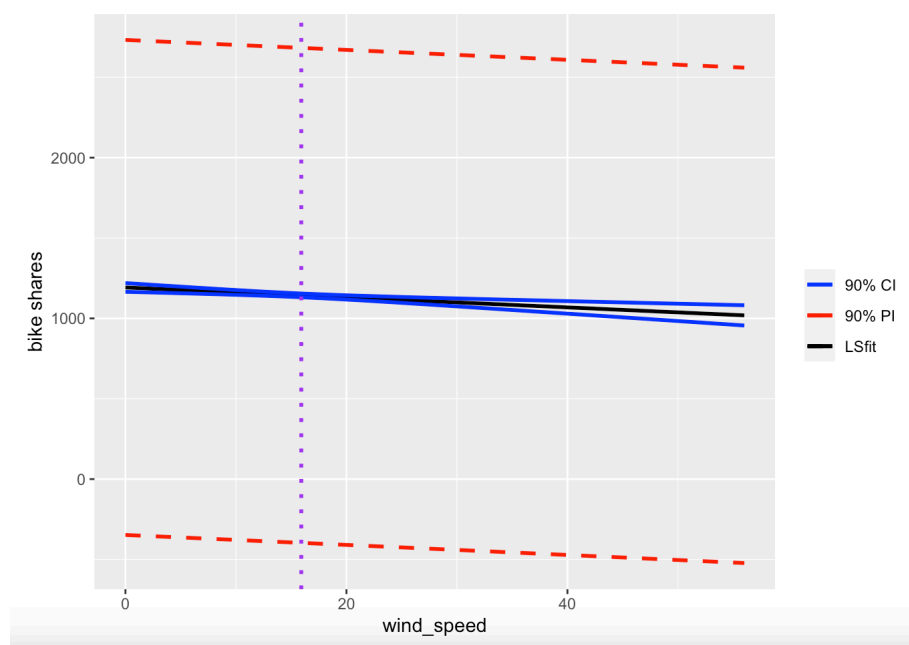


图 19: