

Математические пакеты

Методы обработки данных средней сложности. Регрессия

Сучков Андрей Игоревич

Санкт-Петербургский государственный электротехнический университет «ЛЭТИ»

21 ноября 2020 г.

Дополнение к предыдущей лекции

Диаграммы ядерной оценки функции плотности

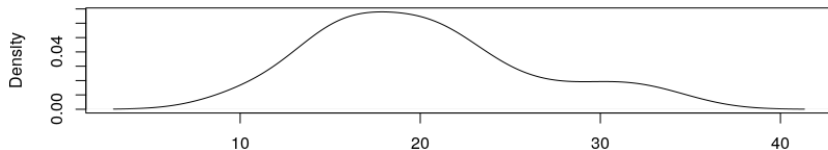
Листинг 1: Пример построения диаграммы ядерной оценки функции плотности

```
1 par(mfrow=c(2,1))
2 d <- density(mtcars$mpg)
3 plot(d)
4 d <- density(mtcars$mpg)
5 plot(d, main="Kernel Density of Miles Per Gallon")
6 polygon(d, col="red", border="blue")
7 rug(mtcars$mpg, col="brown")
```

Дополнение к предыдущей лекции

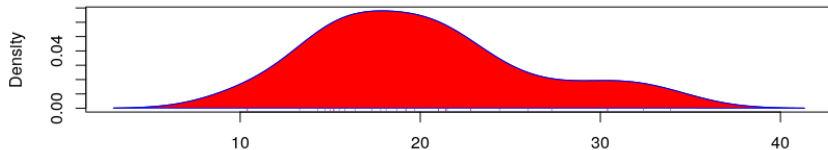
Результат построения диаграммы

density.default(x = mtcars\$mpg)



N = 32 Bandwidth = 2.477

Kernel Density of Miles Per Gallon



N = 32 Bandwidth = 2.477

Дополнение к предыдущей лекции

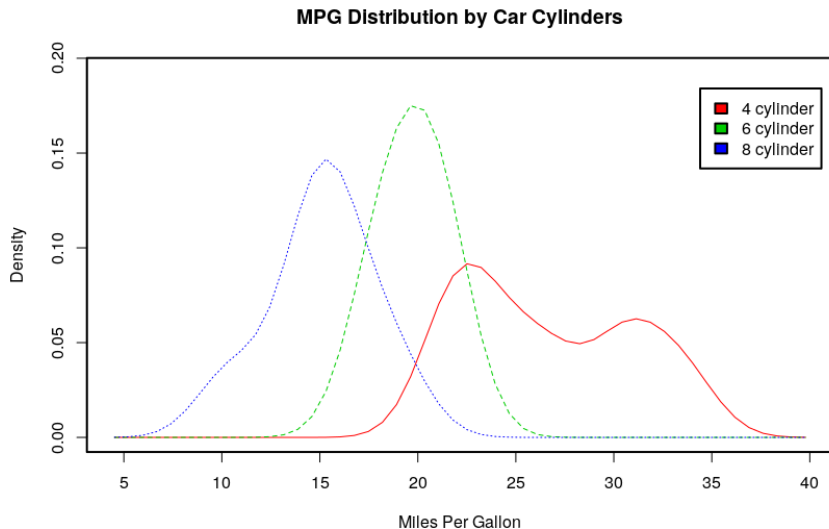
Диаграммы ядерной оценки функции плотности

Листинг 2: Сравнение диаграмм ядерной оценки функции плотности

```
1 par(lwd=2)
2 library(sm)
3 attach(mtcars)
4 cyl.f <- factor(cyl, levels = c(4,6,8),
5                 labels = c("4 cylinder",
6                             "6 cylinder",
7                             "8 cylinder"))
8
9 sm.density.compare(mpg, cyl,
10                   xlab="Miles Per Gallon")
11 title(main="MPG Distribution by Car Cylinders")
12
13 colfill <- c(2:(1+length(levels(cyl.f))))
14 legend(locator(1), levels(cyl.f), fill=colfill)
15 detach(mtcars)
```

Дополнение к предыдущей лекции

Результат построения диаграммы

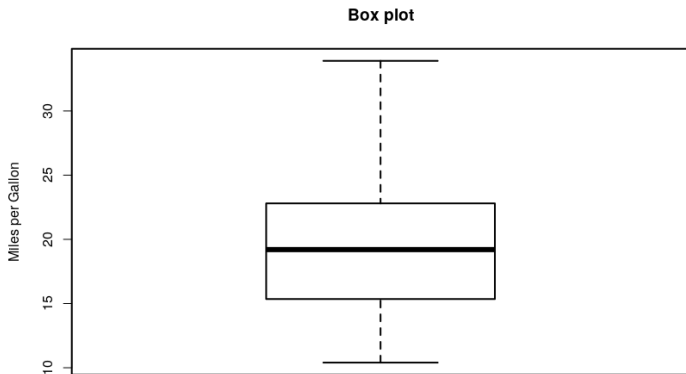


Дополнение к предыдущей лекции

Диаграммы размахов («ящик с усами»)

Построение «ящика с усами» в R

```
boxplot(mtcars$mpg, main="Box plot", ylab="Miles per  
Gallon")
```



Дополнение к предыдущей лекции

Использование диаграмм размахов для сравнения групп между собой

Листинг 3: Сравнение диаграмм размахов

```
1 boxplot(mpg ~ cyl, data=mtcars,
2         main="Car Mileage Data",
3         xlab="Number of Cylinders",
4         ylab="Miles Per Gallon")
```

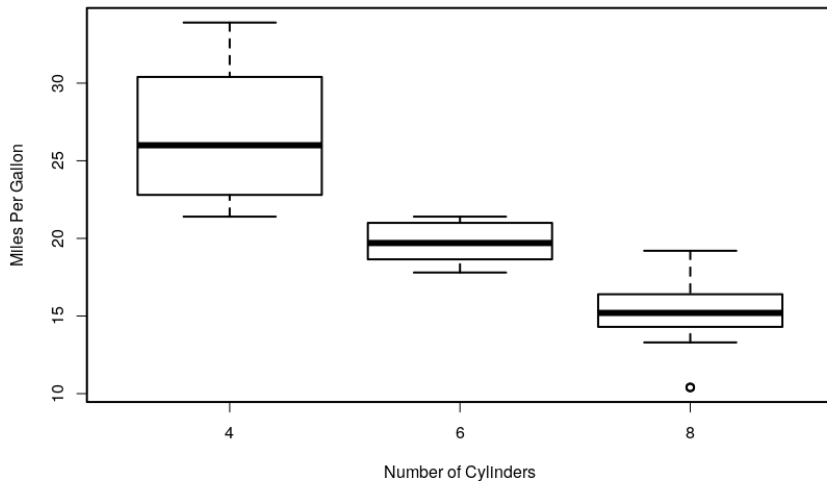
Листинг 4: Построение диаграмм размахов с «насечками»

```
1 boxplot(mpg ~ cyl, data=mtcars,
2         notch=TRUE,
3         varwidth=TRUE,
4         col="red",
5         main="Car Mileage Data",
6         xlab="Number of Cylinders",
7         ylab="Miles Per Gallon")
```

Дополнение к предыдущей лекции

Результат построения диаграммы

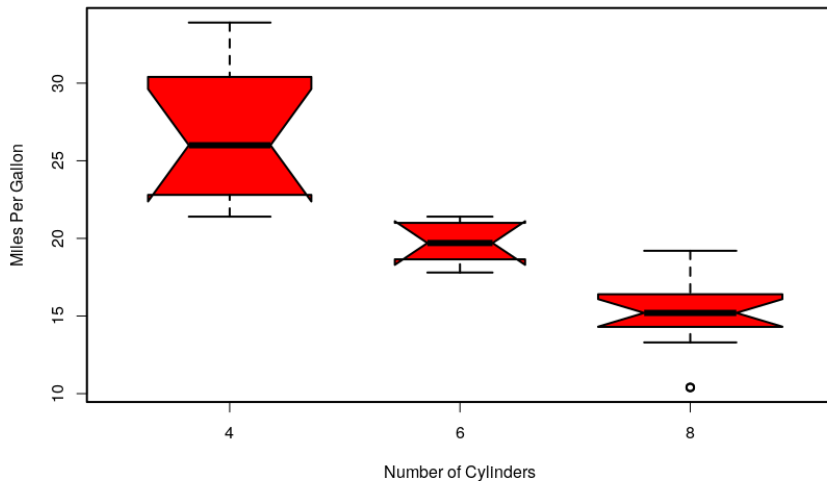
Car Mileage Data



Дополнение к предыдущей лекции

Результат построения диаграммы с «насечками»

Car Mileage Data



Дополнение к предыдущей лекции

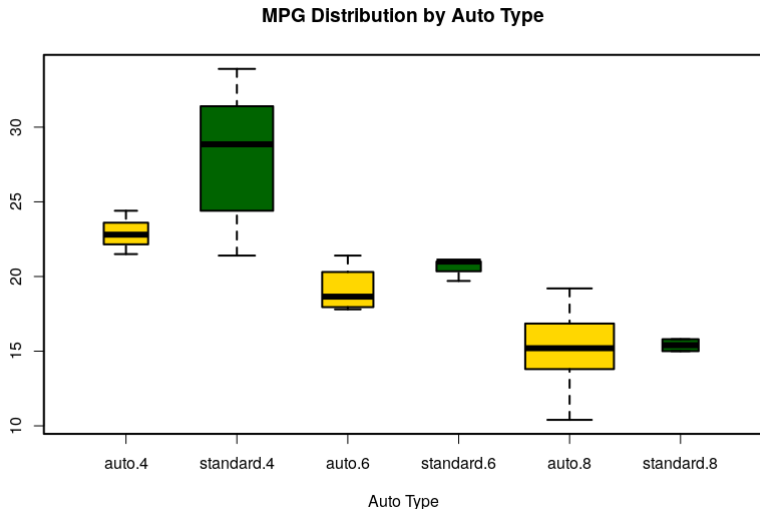
Диаграммы размахов для нескольких группирующих переменных

Листинг 5: Диаграмма размахов для всех сочетаний значений двух факторов

```
1 mtcars$cyl.f <- factor(mtcars$cyl,  
2                       levels=c(4,6,8),  
3                       labels=c("4","6","8"))  
4  
5 mtcars$am.f <- factor(mtcars$am,  
6                      levels=c(0,1),  
7                      labels=c("auto", "standard"))  
8  
9 boxplot(mpg ~ am.f * cyl.f, data=mtcars,  
10        varwidth=TRUE,  
11        col=c("gold","darkgreen"),  
12        main="MPG Distribution by Auto Type",  
13        xlab="Auto Type")
```

Дополнение к предыдущей лекции

Диаграмма размахов, отражающая расход топлива для всех комбинаций числа цилиндров и типа коробки передач



Дополнение к предыдущей лекции

Скрипичные диаграммы

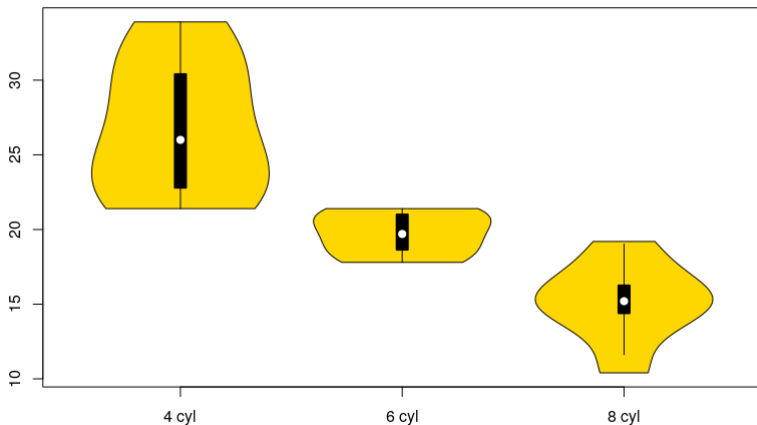
Листинг 6: Построение скрипичных диаграмм

```
1 library(vioplot)
2 x1 <- mtcars$mpg[mtcars$cyl==4]
3 x2 <- mtcars$mpg[mtcars$cyl==6]
4 x3 <- mtcars$mpg[mtcars$cyl==8]
5 violplot(x1, x2, x3,
6           names=c("4 cyl", "6 cyl", "8 cyl"),
7           col="gold")
8 title("Violin Plots of Miles Per Gallon")
```

Дополнение к предыдущей лекции

Скрипичные диаграммы, отражающие расход топлива у автомобилей с разным числом цилиндров

Violin Plots of Miles Per Gallon



Дополнение к предыдущей лекции

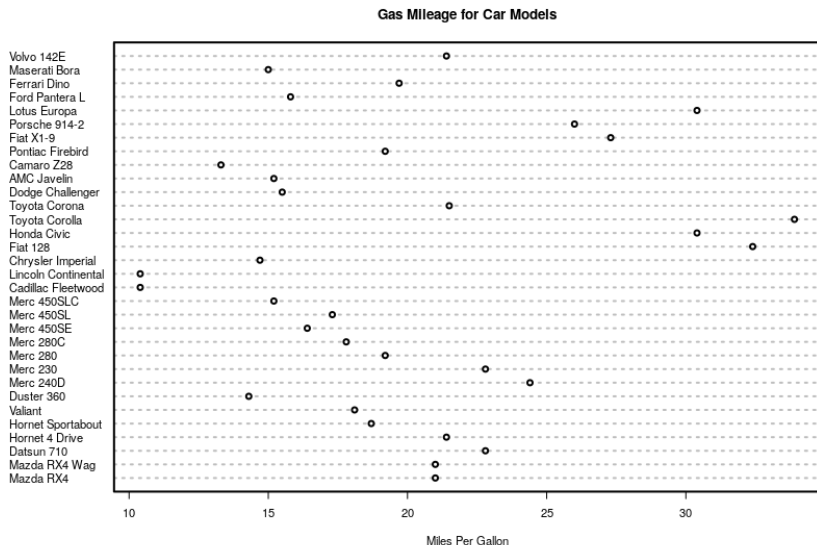
Точечные диаграммы

Листинг 7: Построение точечных диаграмм

```
1 with(mtcars, {  
2   dotchart(mpg, labels=row.names(mtcars),  
3           cex=.7,  
4           main="Gas Mileage for Car Models",  
5           xlab="Miles Per Gallon"))})
```

Дополнение к предыдущей лекции

Точечная диаграмма для расхода топлива у автомобилей разных марок



Дополнение к предыдущей лекции

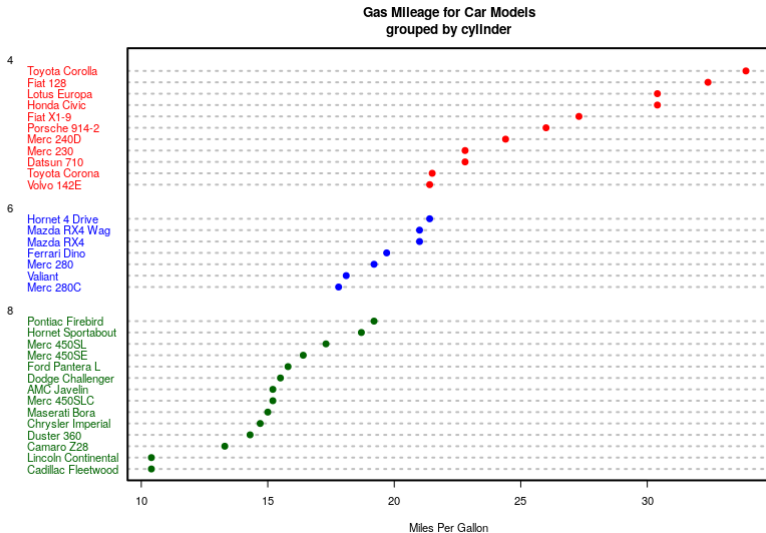
Точечная диаграмма с отсортированными, сгруппированными и раскрашенными значениями

Листинг 8: Построение точечных диаграмм

```
1 x <- mtcars[order(mtcars$mpg),]  
2 x$cyl <- factor(x$cyl)  
3 x$color[x$cyl==4] <- "red"  
4 x$color[x$cyl==6] <- "blue"  
5 x$color[x$cyl==8] <- "darkgreen"  
6 dotchart(x$mpg,  
7           labels = row.names(x),  
8           cex=.7,  
9           groups = x$cyl,  
10          gcolor = "black",  
11          color = x$color,  
12          pch=19,  
13          main = "Gas Mileage for Car Models\ngrouped by  
14             cylinder",  
             xlab = "Miles Per Gallon")
```


Дополнение к предыдущей лекции

Точечная диаграмма для расхода топлива у автомобилей, сгруппированных по числу цилиндров



Вид модели МНК-регрессии

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \dots + \hat{\beta}_k X_{ki} \quad i = 1..n$$

- \hat{Y}_i – предсказанное значение зависимой переменной для i -го наблюдения (а именно оценка среднего значения распределения Y по набору независимых переменных)
- X_{ki} – значение k -ой независимой переменной для i -го наблюдения
- $\hat{\beta}_0$ – свободный член уравнения (предсказанное значение Y при нулевом значении всех независимых переменных)
- $\hat{\beta}_k$ – регрессионный коэффициент для k -ой независимой переменной (угол наклона для прямой, которая отражает изменение Y при изменении X на одну единицу измерения)

Цель

Выбрать такие параметры модели (свободный член и регрессионные коэффициенты), которые позволят минимизировать различия между реальными и предсказанными значениями зависимой переменной. То есть мы выбираем такие параметры модели, чтобы сумма квадратов остатков была минимальной:

$$\sum_1^n (Y_i - \hat{Y}_i)^2 = \sum_1^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \dots - \hat{\beta}_k X_{ki})^2 = \sum_1^n \varepsilon^2$$

Требования:

- нормальность
- независимость
- линейность
- гомоскедастичность

Регрессионный анализ

Подгонка регрессионных моделей при помощи команды `lm()`

Формат функции `lm()`

```
myfit <- lm(formula, data)
```

Формула обычно записывается в таком виде:

$$Y \sim X_1 + X_2 + \dots + X_k$$

Символы, которые часто используются в формулах R:

- `~` – Отделяет зависимые переменные (слева) от независимых (справа)
- `+` – Разделяет независимые переменные
- `:` – Обозначает взаимодействие между независимыми переменными
- `*` – Краткое обозначение для всех возможных взаимодействий
- `^` – Обозначает взаимодействия до определенного порядка

Регрессионный анализ

Подгонка регрессионных моделей при помощи команды `lm()`

Символы, которые часто используются в формулах R:

- `.` – Символ-заполнитель для всех переменных в таблице данных, кроме зависимой
- `-` – Знак минуса удаляет переменную из уравнения
- `-1` – Подавляет свободный член уравнения
- `I()` – Элемент в скобках интерпретируется как арифметическое выражение
- `function` – В формулах можно использовать математические функции

Листинг 9: Простая линейная регрессия

```
1 fit <- lm(weight ~ height, data=women)
2 summary(fit)
3
4 women$weight
5
6 fitted(fit)
7
8 residuals(fit)
9
10 plot(women$height, women$weight,
11       xlab="Height (in inches)",
12       ylab="Weight (in pounds)")
13 abline(fit)
```

Регрессионный анализ

Результат выполнения команды `summary(fit)`

```
Call:
lm(formula = weight ~ height, data = women)

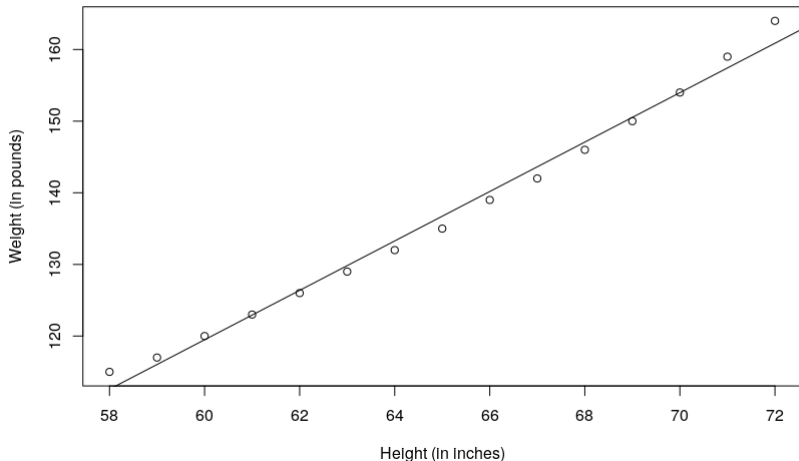
Residuals:
    Min       1Q   Median       3Q      Max
-1.7333 -1.1333 -0.3833  0.7417  3.1167

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -87.51667    5.93694  -14.74 1.71e-09 ***
height        3.45000    0.09114   37.85 1.09e-14 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.525 on 13 degrees of freedom
Multiple R-squared:  0.991,    Adjusted R-squared:  0.9903
F-statistic: 1433 on 1 and 13 DF,  p-value: 1.091e-14
```

Регрессионный анализ

Диаграмма рассеяния с регрессионной прямой для значений веса, предсказанных по значениям роста



Листинг 10: Полиномиальная регрессия

```
1 fit2 <- lm(weight ~ height + I(height^2), data=women)
2 summary(fit2)
3
4 plot(women$height, women$weight,
5       xlab="Height (in inches)",
6       ylab="Weight (in lbs)")
7 lines(women$height, fitted(fit2))
```

Регрессионный анализ

Результат выполнения команды `summary(fit2)`

```
Call:
lm(formula = weight ~ height + I(height^2), data = women)

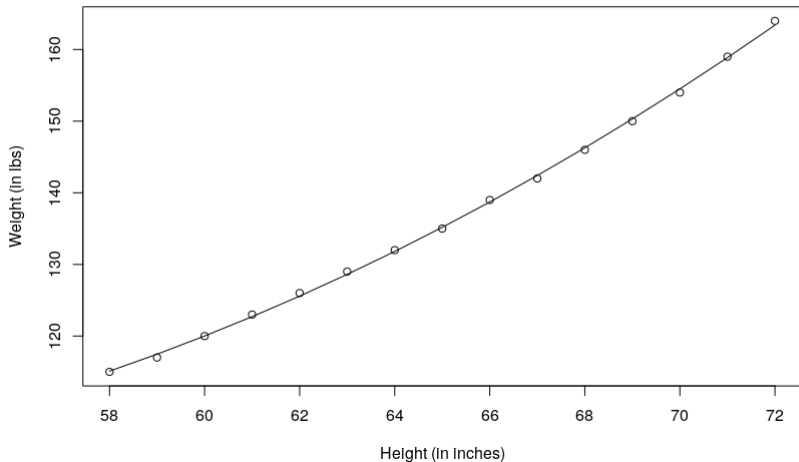
Residuals:
      Min       1Q   Median       3Q      Max
-0.50941 -0.29611 -0.00941  0.28615  0.59706

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  261.87818   25.19677   10.393 2.36e-07 ***
height       -7.34832    0.77769   -9.449 6.58e-07 ***
I(height^2)   0.08306    0.00598   13.891 9.32e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3841 on 12 degrees of freedom
Multiple R-squared:  0.9995,    Adjusted R-squared:  0.9994
F-statistic: 1.139e+04 on 2 and 12 DF,  p-value: < 2.2e-16
```

Регрессионный анализ

Квадратичная регрессия для предсказаний значений веса по значениям роста



Диагностика регрессионных моделей

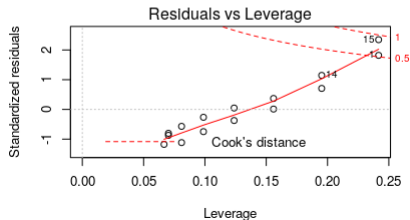
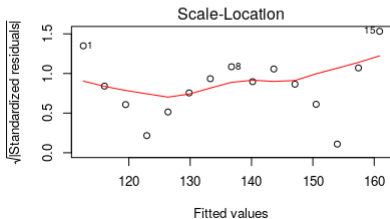
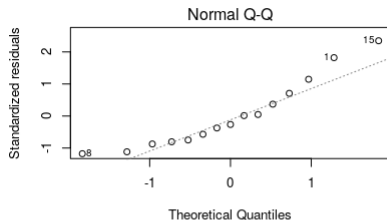
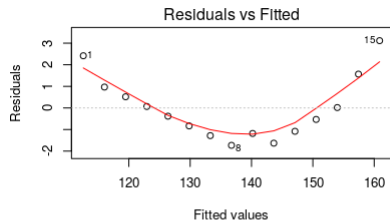
Стандартный подход

Листинг 11: Полиномиальная регрессия

```
1 fit <- lm(weight ~ height, data=women)
2 par(mfrow=c(2,2))
3 plot(fit)
4
5 fit2 <- lm(weight ~ height + I(height^2), data=women)
6 par(mfrow=c(2,2))
7 plot(fit2)
8
9 states <- as.data.frame(state.x77[,c("Murder", "Population",
10                                     "Illiteracy", "Income",
11                                     "Frost")])
12 fit3<-lm(Murder ~ Population + Illiteracy + Income + Frost,
13          data=states)
14 par(mfrow=c(2,2))
15 plot(fit3)
```

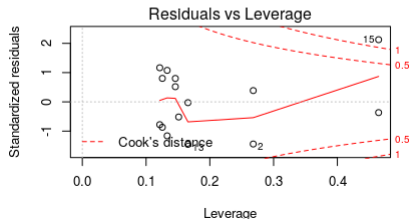
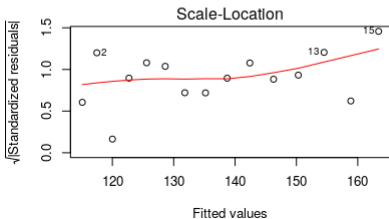
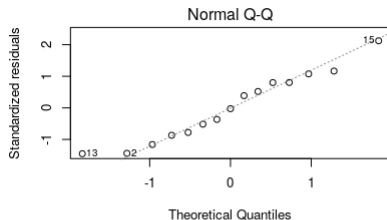
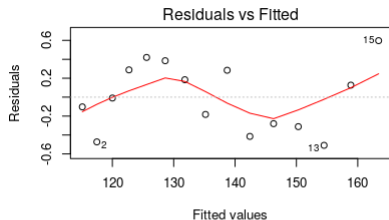
Диагностика регрессионных моделей

Диагностические диаграммы для регрессии веса по росту



Диагностика регрессионных моделей

Диагностические диаграммы для регрессии веса по росту и по росту, возведенному в квадрат



Диагностика регрессионных моделей

Диагностические диаграммы для регрессии уровня преступности по характеристикам штатов

