

Линейная регрессия

Модель

Данные: пары чисел $\{x_i, y_i\}$

Гипотеза: имеется линейная статистическая зависимость между переменными X и Y

$$Y = a + b \cdot X + \varepsilon$$

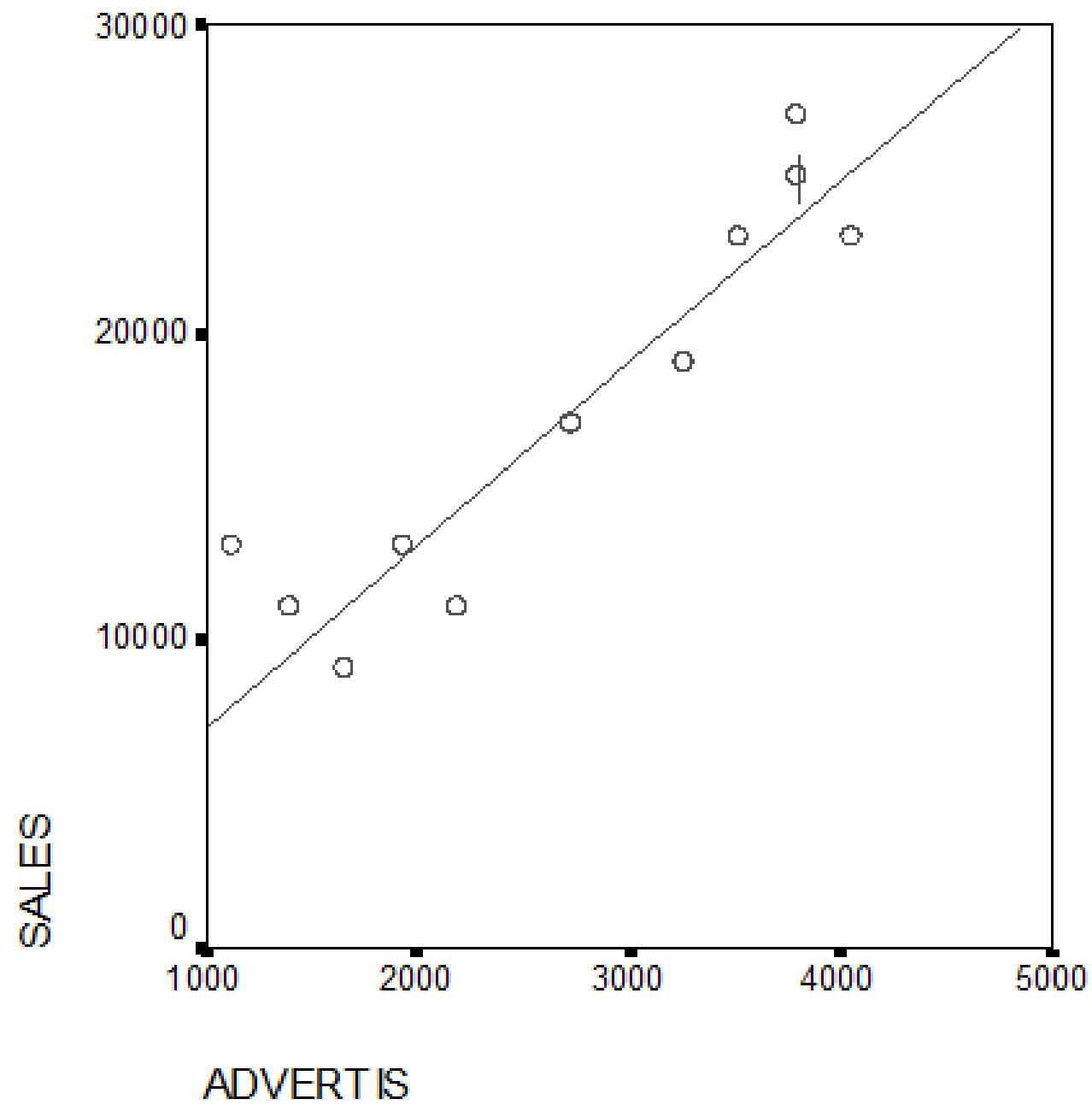
Надо найти

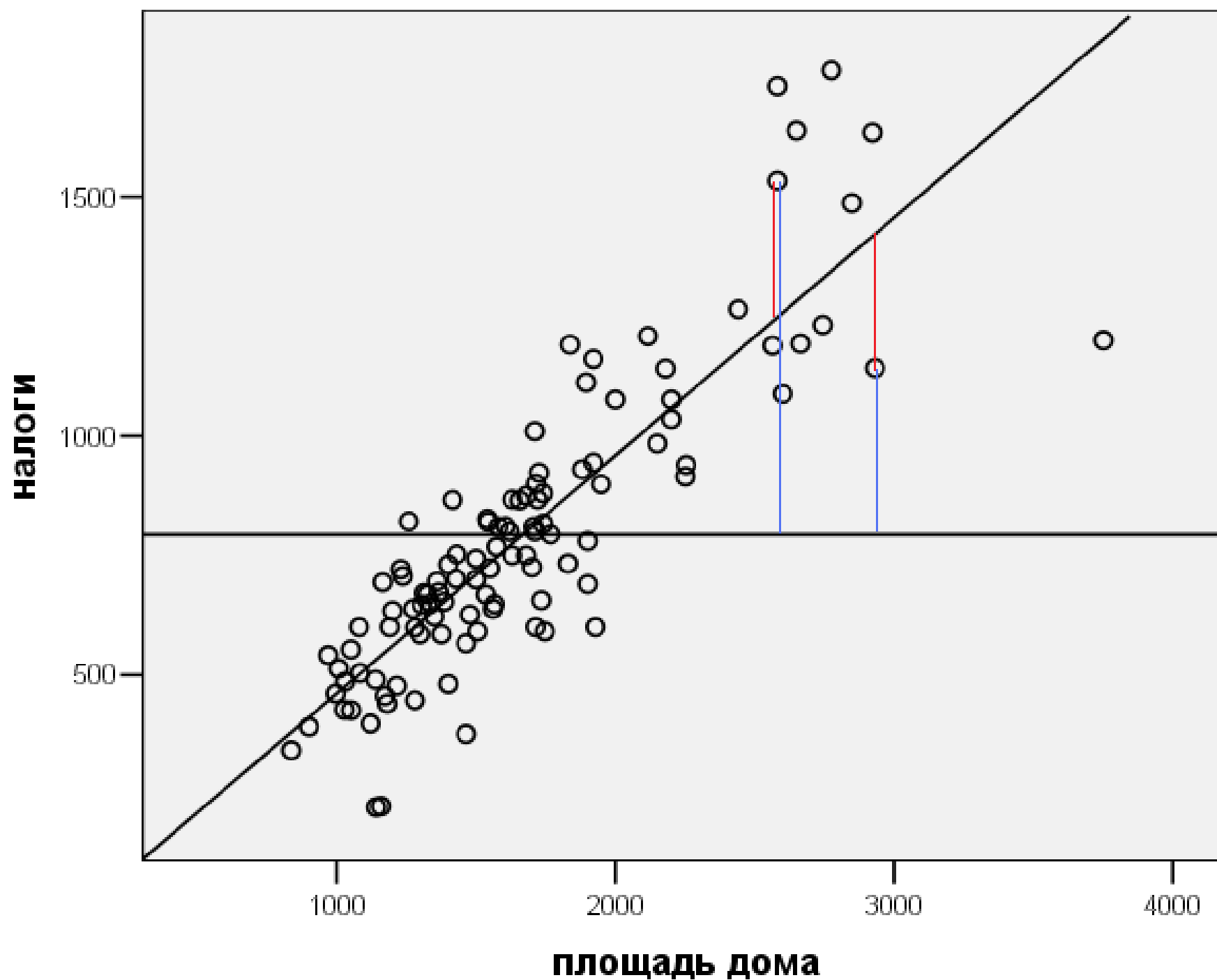
оценки коэффициентов a и b уравнения
регрессии:

$$Y = a + b \cdot X$$

Геометрический взгляд

Уравнение определяет прямую, наиболее близко проходящую ко всем точкам с координатами x_i, y_i





Алгебраический взгляд

Значения a и b находятся по методу наименьших квадратов, т.е. так, чтобы минимизировать величину

$$SSE = \frac{1}{n} \sum_{i=1}^n (y_i - (a + b \cdot x_i))^2$$

Терминология

- Y — отклик / зависимая переменная
- X — предиктор / независимая переменная
- $\varepsilon_i = y_i - (a + b \cdot x_i)$ - ошибка / невязка /
-

Измеряем качество модели

Регрессионная модель

Сумма квадратов отклонений для
регрессионной модели

Зависит от единиц измерения

$$SSE = \frac{1}{n} \sum_{i=1}^n (y_i - (a + b \cdot x_i))^2$$

Измеряем качество модели

Базовая модель

Сумма квадратов отклонений для базовой модели

Зависит от единиц измерения

$$SSE = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)^2$$

Коэффициент детерминации

$$R^2 = 1 - \frac{\frac{1}{n} \sum_{i=1}^n (y_i - (a + b \cdot x_i))^2}{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)^2}$$

Коэффициент детерминации

указывает, какой процент вариации у объясняется влиянием предикторов.

Недостатки множественного коэффициента детерминации

Пример «Квартет Ансcombe»

Anscombe, F. J. (1973). "Graphs in Statistical Analysis". American Statistician 27 (1): 17–21.

Квартет Анскомба: 4 набора данных

Характеристика	Значение
Mean of x in each case	9
Sample variance of x in each case	11
Mean of y in each case	7.50
Sample variance of y in each case	4.122 or 4.127
Correlation between x and y in each case	0.816
Linear regression line in each case	$y = 3.00 + 0.500x$

Квартет Анскомба в R

- `help(anscombe)`

```
ans <- anscombe
```

```
res1 <- lm(ans$y1~ans$x1)
```

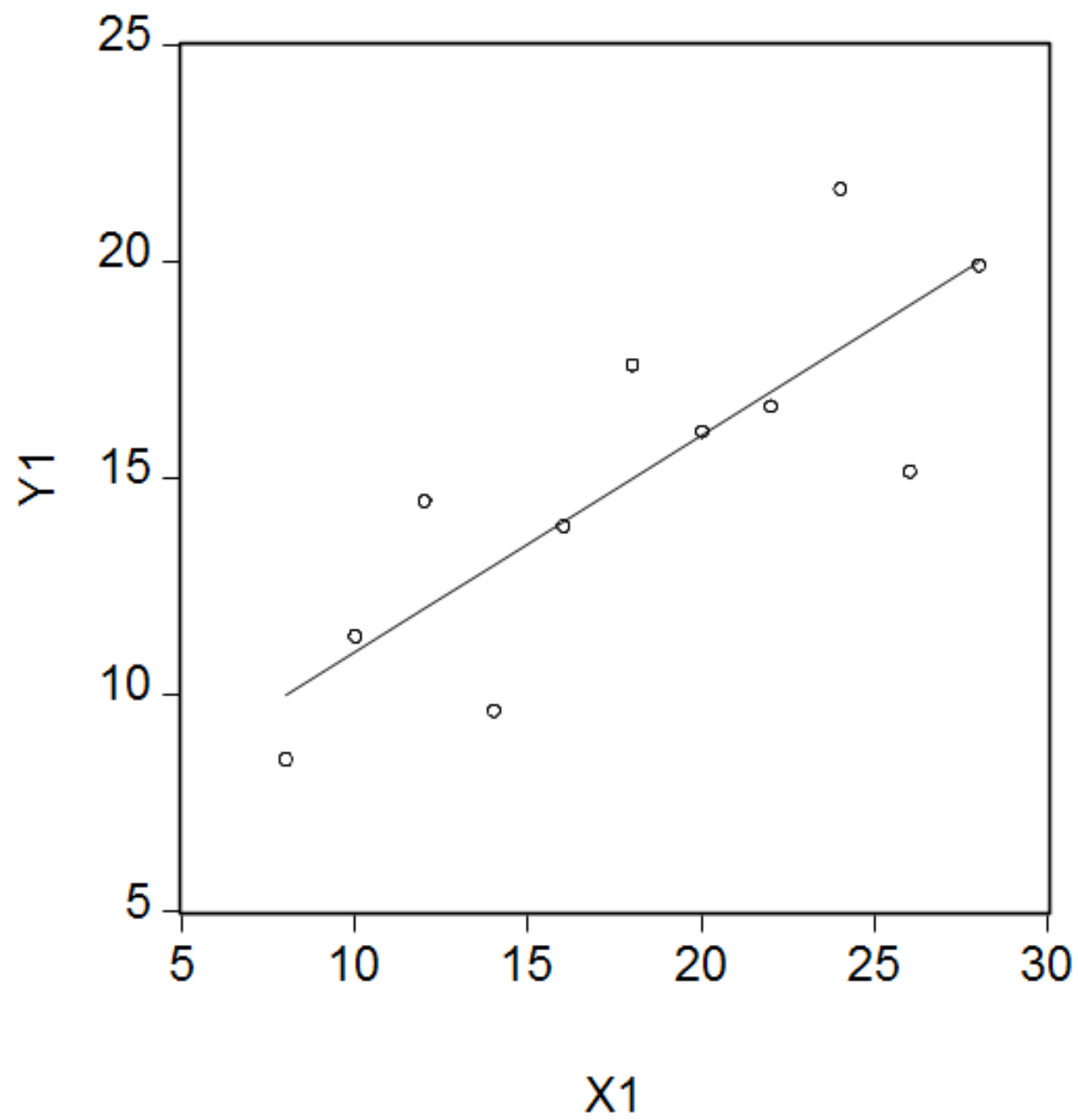
```
plot(ans$x1,ans$y1)
```

```
abline(res1)
```

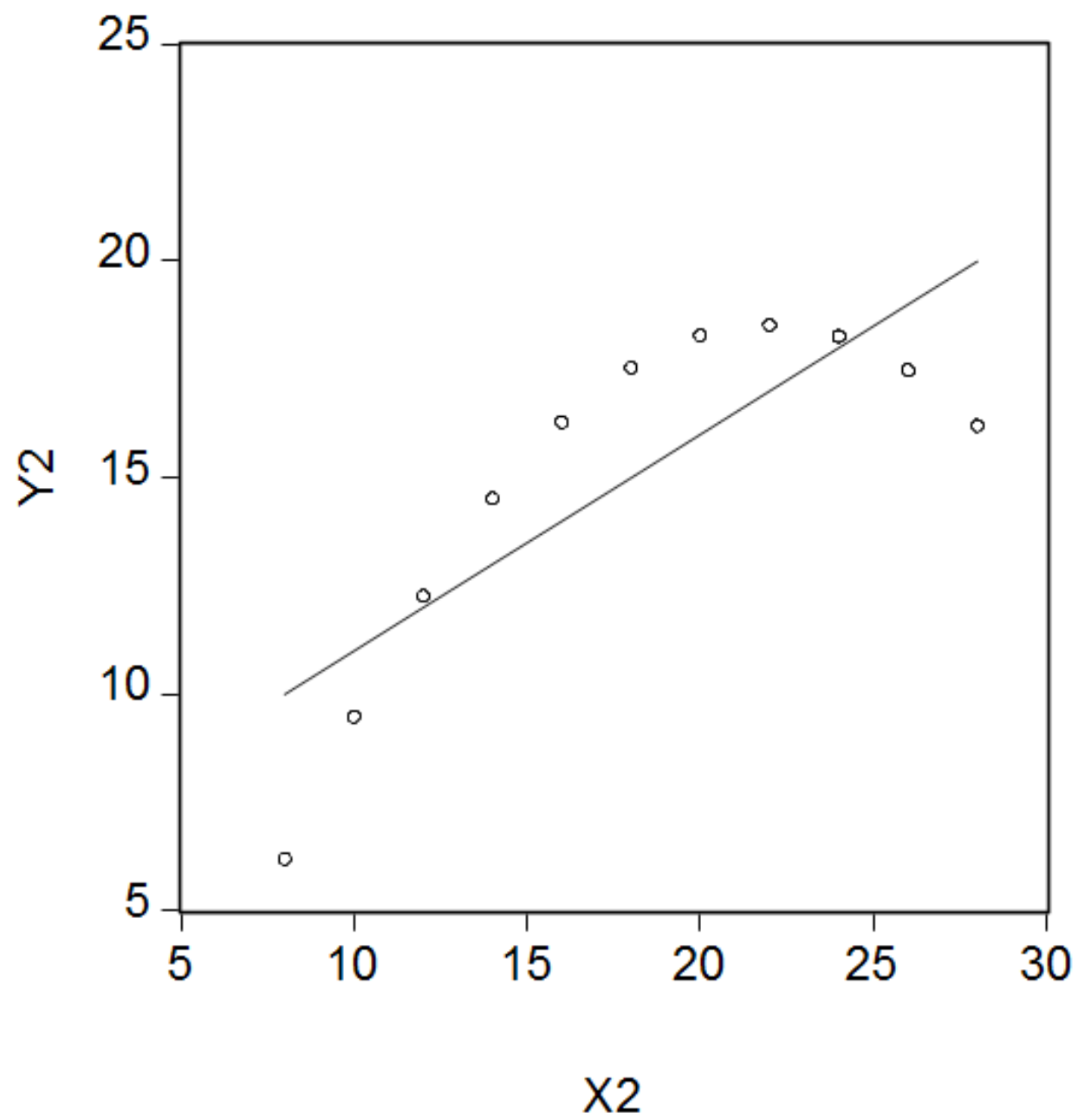
-

-

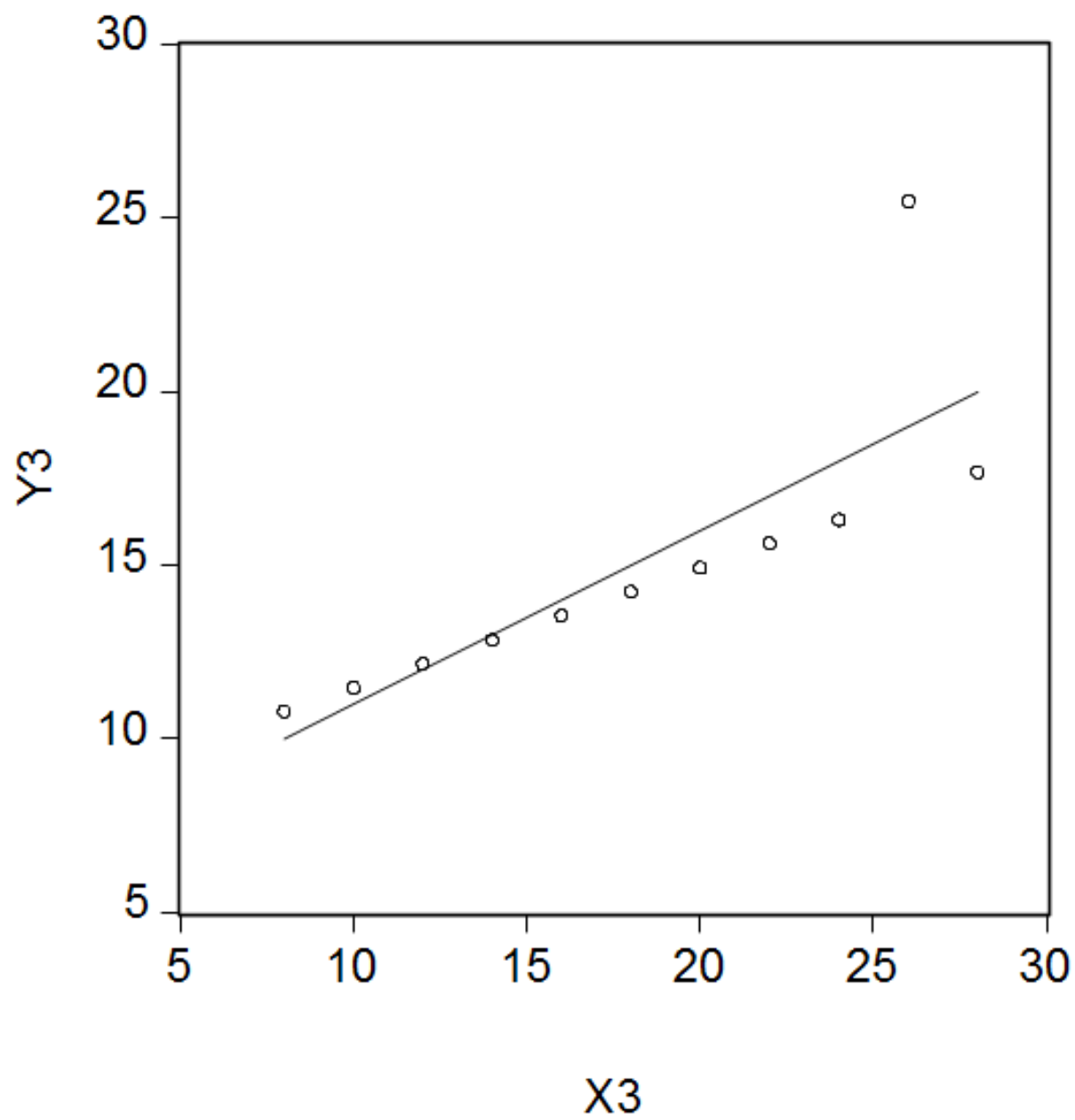
Y1 vs. X1



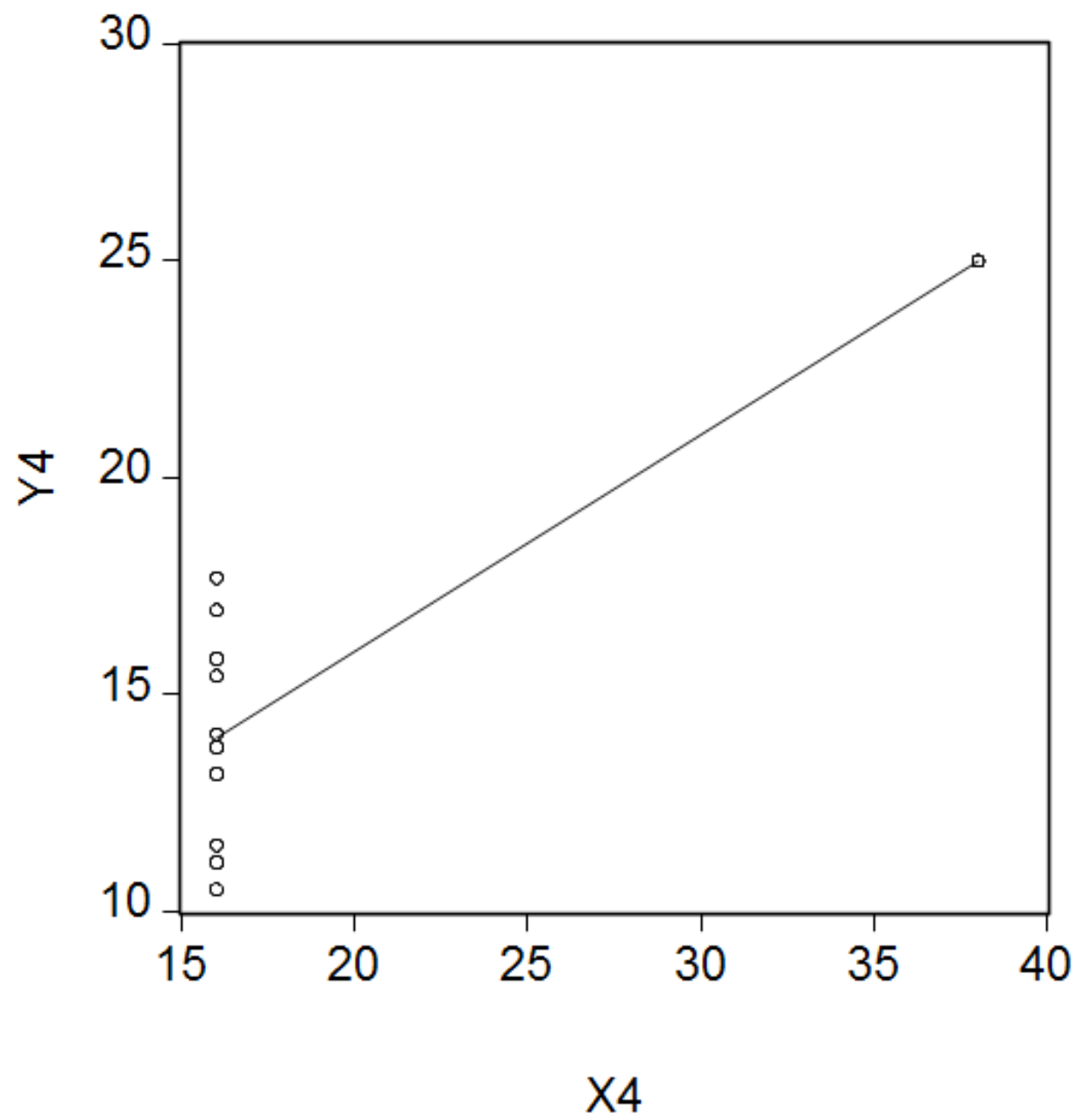
Y2 vs. X2



Y3 vs. X3



Y4 vs. X4



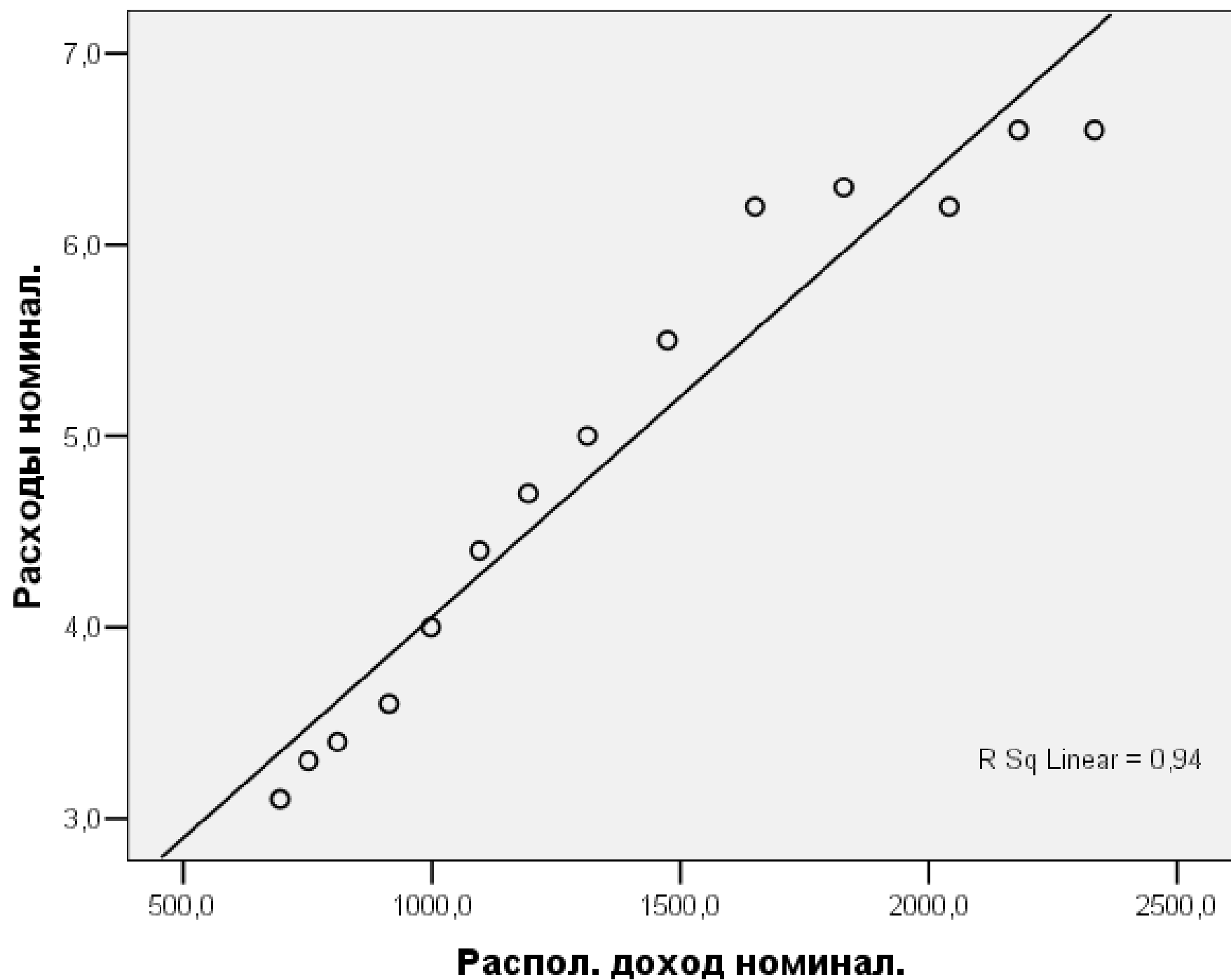
Влияние инфляции на модель

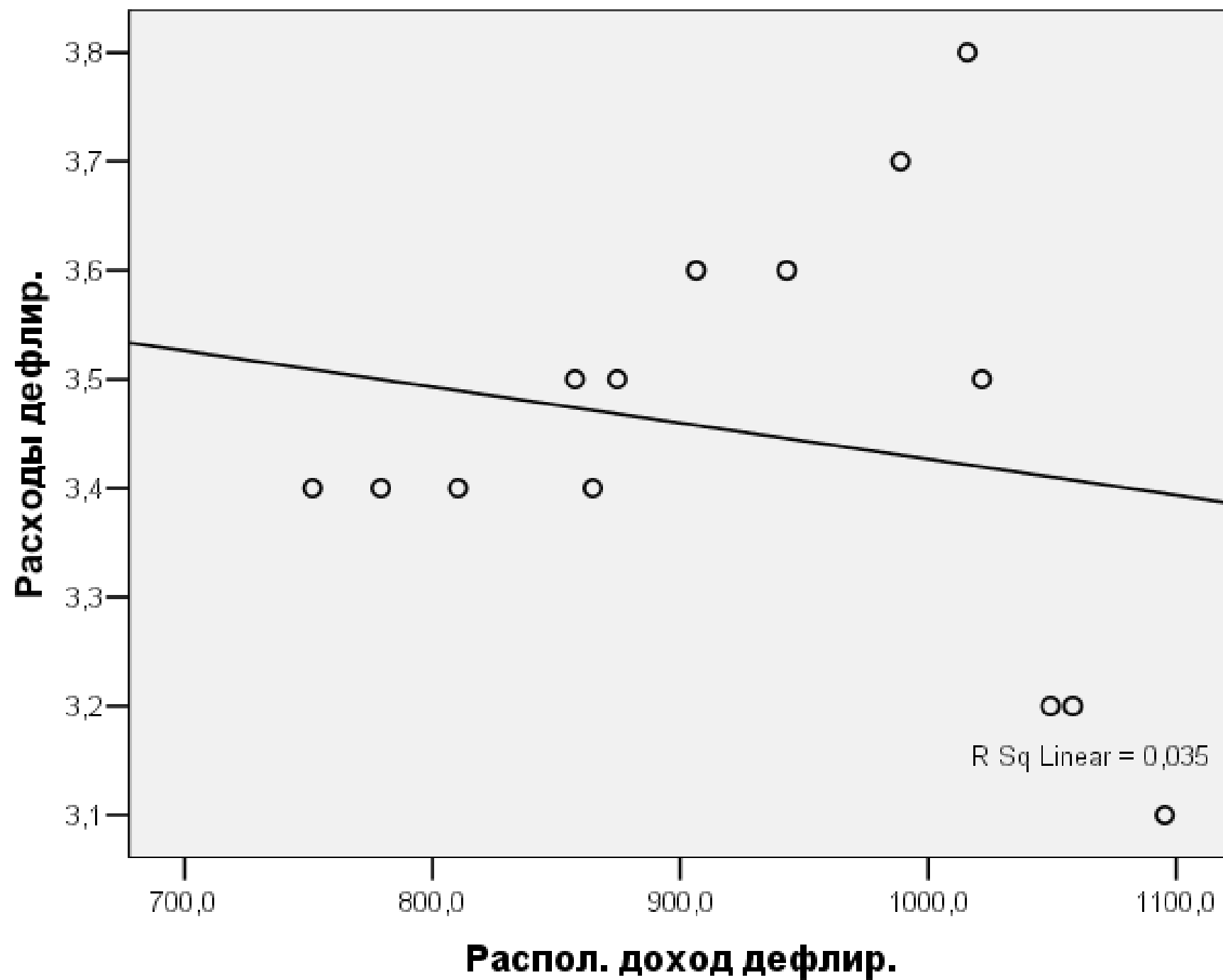
Пример

совокупный располагаемый доход и
совокупные личные расходы на местный транспорт
в США за период с 1970 по 1983 год.

Данные представлены как в текущих долларах США, так и в долларах 1972 года — пересчет к последним выполнен с учетом динамики индекса потребительских цен в указанном периоде.

(Уровень цен в 1972 г. принят за 100%.)





Коллинеарность-1

Допустим, что задана зависимость

$$Z = 6X + 8Y + 5$$

Одновременно известно, что

$$Y = 2X$$

Коллинеарность-2

Если верно $Z = 6X + 8Y + 5$,

То верно

$$Z = 2X + 10Y + 5$$

$$Z = 22X + 0 \cdot Y + 5$$

$$Z = 0 \cdot X + 11 \cdot Y + 5$$

И так далее...

(На самом деле все еще хуже...)

tolerance
пороговые значения 0.1 или 0.2

$$tolerance = 1 - R_j^2$$

VIF variance inflation factor

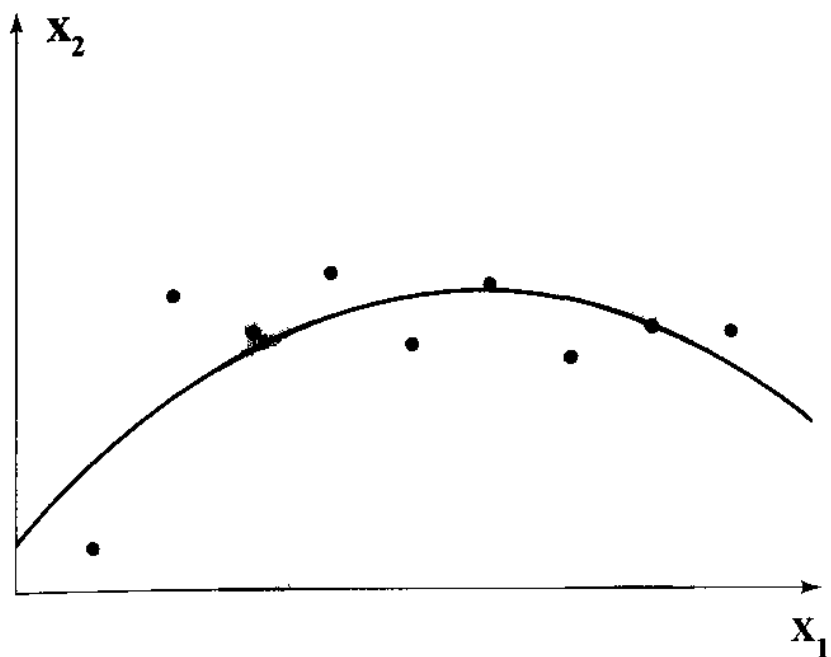
пороговые значения 5 или 10

$$VIF = \frac{1}{tolerance}$$

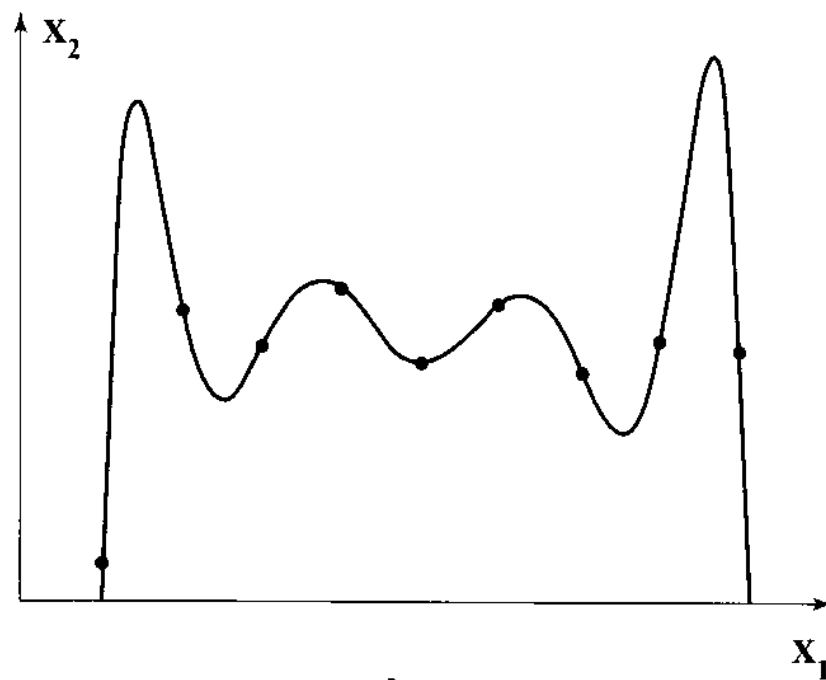
Вычисления VIF в R

- `library(car)`
- `# пример Albuquerque Home Prices`
- `vif(fit)` `# variance inflation factors`
- `vif(itog1)`
- | <code>x\$SQFT</code> | <code>x\$AGE</code> | <code>x\$FEATS</code> | <code>x\$NE</code> | <code>x\$CUST</code> | <code>x\$COR</code> | <code>x\$TAX</code> |
|----------------------|---------------------|-----------------------|--------------------|----------------------|---------------------|---------------------|
| 6.197 | 1.692 | 1.459 | 1.374 | 1.385 | 1.108 | 6.476 |

Переобучение



a



b

Non-independence of Errors

-
-
- # Test for Autocorrelated Errors
- `durbinWatsonTest(fit)`