



# Прикладная статистика и анализ данных

Съезд II



# Анализ остатков

# Остатки

В качестве оценки ошибки  $\varepsilon_i$  рассмотрим остатки  $e_i = Y_i - \hat{Y}_i$

## Проверка свойств

### Нормальность

$$H_0: e_i \sim \mathcal{N}$$



Критерий Шапиро-Уилка и др.

### Несмещенность

$$H_0: Ee_i = 0$$



Критерии монотонного отнош. правд.

В непарам. случае позже

### Гомоскедастичность

$$H_0: D\varepsilon = \sigma^2 I_n$$



Тут не все так просто...



## Остатки

$D\varepsilon = \sigma^2 I_n$  — гомоскедастичность. Обратное — гетероскедастичность.  
В качестве оценки ошибки  $\varepsilon_i$  рассмотрим остатки  $e_i = Y_i - \hat{Y}_i$

**Проблема:**  $De_i \neq \sigma^2$  при гомоскедастичности.

$$e = Y - \hat{Y} = (I_n - H)Y, \quad \text{где } H = X(X^T X)^{-1}X^T$$

$$De = (I_n - H)DY(I_n - H)^T = \sigma^2(I_n - H)(I_n - H)^T = \sigma^2(I_n - H)$$

Проверять на однородность дисп. нужно **поправленные остатки**:

$$\hat{e}_i = \frac{e_i}{\sqrt{De_i}} = \frac{e_i}{\sqrt{\frac{RSS}{n-d}(1 - h_{ii})}} \text{ — студентизированные остатки}$$

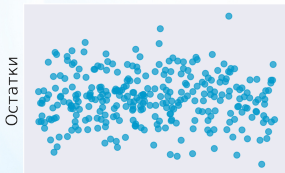
Т.к.  $\text{tr } H = d$  [упр.], то при  $d \ll n$ :  $h_{ii} \approx 0$ . Тогда

$$\hat{e}_i = \frac{e_i}{\sqrt{RSS/(n-d)}} \text{ — стандартизированные остатки}$$



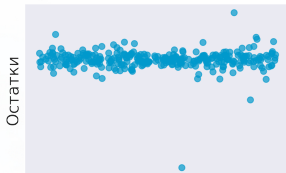
# Визуальный анализ

Строятся графики зависимости  $\hat{e}_i$  от  $y, x, i$



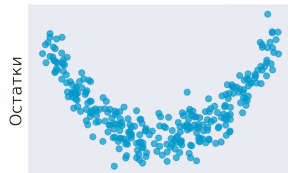
Признак

Все хорошо



Предсказание

Есть выбросы



Остатки

Признак

Нужно добавить  $x^2$



Признак

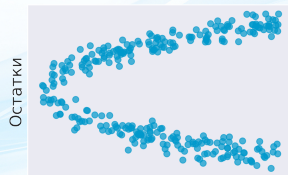
Гетероскедастичность



Остатки

Номер наблюдения

Тренд



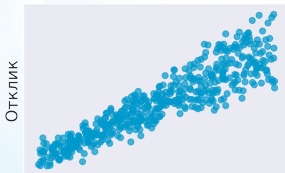
Остатки

Признак

Неправильная модель

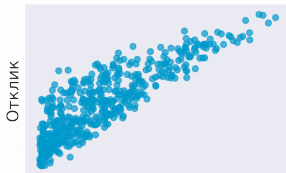
# Визуальный анализ

Что будет если строить графики зависимостей остатков от признаков:



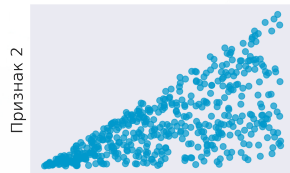
Признак 1

Гетероскедастичность?



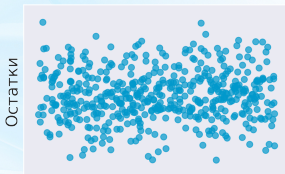
Признак 2

Гетероскедастичность?



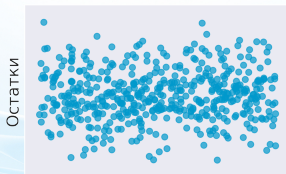
Признак 1

Признаки зависимы

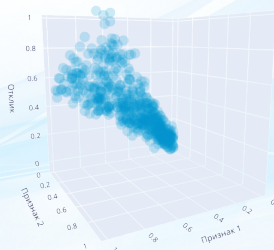


Признак 1

Нет, все хорошо!



Признак 2





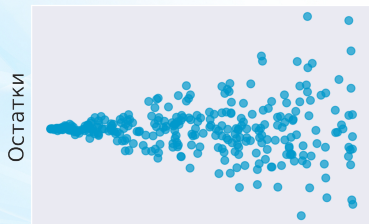
# Критерии проверки на гомоскедастичность

$$H_0: D\varepsilon = \sigma^2 I_n$$

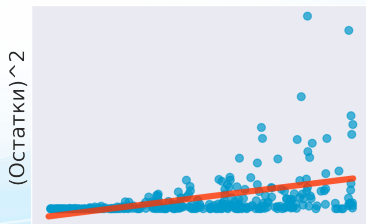
## Критерий Бройша-Пагана

$R_{\hat{e}^2}^2$  — коэф. детерминации при регрессии  $\hat{e}^2 \sim X$

$nR_{\hat{e}^2}^2 \sim \chi_d^2$  — при справедливости  $H_0$



Признак



Признак



# Критерии проверки на гомоскедастичность

$$H_0: D\varepsilon = \sigma^2 I_n$$

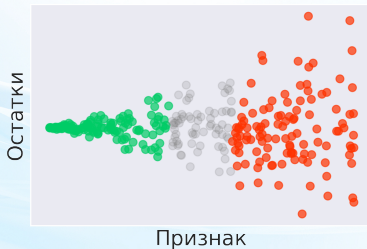
## Критерий Голдфелда-Квандта

Упорядочим наблюдения по предполагаем. возрастанию дисперсий.

$RSS_1$  — регрессия по первым  $\frac{n-r}{2}$  наблюдений,  $r > 0$

$RSS_2$  — регрессия по последним  $\frac{n-r}{2}$  наблюдений

$$\frac{RSS_2}{RSS_1} \sim F_{\frac{n-r}{2}-d, \frac{n-r}{2}-d} \quad \text{при } H_0$$





## Что делать при гетероскедастичности?

- ▶ Если нужна только оценка  $\theta$  — ничего;
- ▶ Если есть предположения о природе гетероскедастичности, взвесить наблюдения:

$$Y_i / \hat{\sigma}_i = (x_i / \hat{\sigma}_i)^T \theta + \varepsilon_i,$$

где  $\hat{\sigma}_i$  — предполагаемая дисперсия при  $i$ -м измерении;

- ▶ Преобразование признаков и отклика, напр., Бокса-Кокса:

$$Z_i = \begin{cases} \ln Y_i, & \lambda = 0 \\ (Y_i^\lambda - 1)/\lambda, & \lambda \neq 0 \end{cases}$$

Величина  $\lambda$  подбирается по графику зависимости  $RSS(\lambda)$  от  $\lambda$

- ▶ Использовать специальные оценки дисперсии, устойчивые к гетероскедастичности.

## Устойчивые оценки дисперсии

Пусть  $E\varepsilon = 0$  и  $D\varepsilon = V$ .

Тогда  $D\hat{\theta} = (X^T X)^{-1} X^T V X (X^T X)^{-1}$ .

1.  $V = \sigma^2 I_n$  — гомоскедастичность:

$D\hat{\theta} = \sigma^2 (X^T X)^{-1}$  — дисперсия оценки коэффициентов;

$\widehat{D\hat{\theta}} = \hat{\sigma}^2 (X^T X)^{-1}$  — оценка дисперсии оценки коэффициентов;

2.  $V = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$  — отсутствие автокорреляций:

$D\hat{\theta} = (X^T X)^{-1} X^T \cdot \text{diag}(\sigma_1^2, \dots, \sigma_n^2) \cdot X (X^T X)^{-1}$  — д.о.к.;

$\widehat{D\hat{\theta}} = (X^T X)^{-1} X^T \cdot \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_n^2) \cdot X (X^T X)^{-1}$  — о.д.о.к..

3. Наличие автокорреляций — см. временные ряды.

## Оценки Уайта

Если автокорреляции отсутствуют, используются оценка Уайта  
White's heteroscedasticity-consistent estimator (HCE):

$$\widehat{D\theta} = (X^T X)^{-1} X^T \cdot \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_n^2) \cdot X (X^T X)^{-1}$$

Варианты определения  $\hat{\sigma}_i^2$ :

1. HC0:  $\hat{e}_i^2$  — оценка Уайта
2. Модификации МакКиннона-Уайта:

$$\text{HC1: } \frac{n}{n-d} \hat{e}_i^2, \quad \text{HC2: } \frac{\hat{e}_i^2}{1-h_{ii}}, \quad \text{HC3: } \frac{\hat{e}_i^2}{(1-h_{ii})^2}$$

(точнее оценивают при малых выборках)



# Асимптотическая нормальность при гетероскедаст.

Если автокорреляции отсутствуют, то

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} \mathcal{N}(0, \Sigma),$$

НСЕ дает состоятельную оценку на матрицу  $\Sigma$ :

$$n\widehat{D\hat{\theta}} \xrightarrow{P} \Sigma$$

Данный факт позволяет строить асимптотические дов. интервалы и критерий Вальда для проверки линейных гипотез  $H_0: T\theta = \tau$ .



**ВСЁ!**