



Статистика

DS-поток

Лекция 13



7.2 Анализ остатков



Остатки

В качестве оценки шума ε_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

$$H_0: E e_i = 0$$



Критерии монотонного отнош. правд.
В непарам. случае позже

Гомоскедастичность

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



Тут не все так просто...
Часто это свойство более критично



Остатки

$D\varepsilon = \sigma^2 I_n$ — гомоскедастичность. Обратное — гетероскедастичность.

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

Проблема: $D e_i \neq \sigma^2$ при гомоскедастичности.

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

$$D e = (I_n - H)D Y(I_n - H)^T = \sigma^2(I_n - H)(I_n - H)^T = \sigma^2(I_n - H)$$

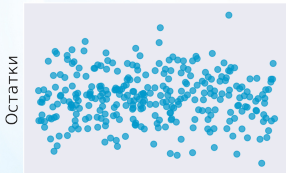
Проверять на однородность дисп. нужно **поправленные остатки**:

$$\hat{e}_i = \frac{e_i}{\sqrt{D e_i}} = \frac{e_i}{\sqrt{\frac{\|Y - X\hat{\theta}\|^2}{n-d}(1 - H_{ii})}} \quad \text{— студентизированные остатки}$$



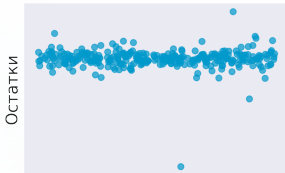
Визуальный анализ

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



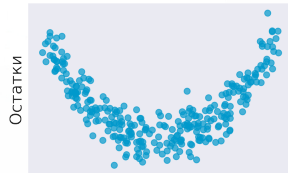
Признак

Все хорошо



Предсказание

Есть выбросы



Признак

Нужно добавить x^2



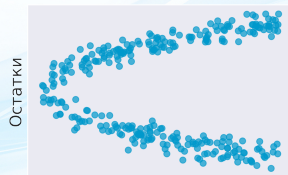
Признак

Гетероскедастичность



Номер наблюдения

Тренд



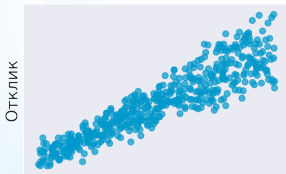
Признак

Неправильная модель



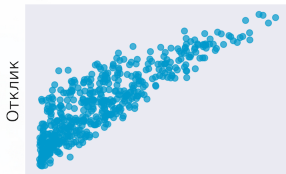
Визуальный анализ

Что будет если строить графики зависимостей таргета от признаков:



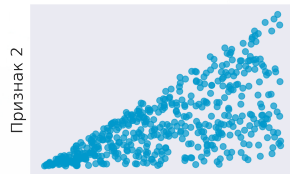
Признак 1

Гетероскедастичность?



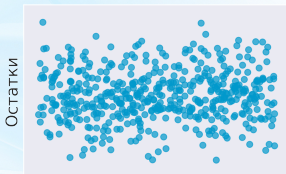
Признак 2

Гетероскедастичность?



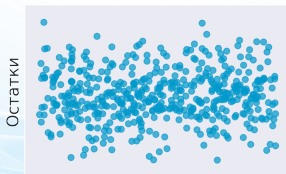
Признак 1

Признаки зависимы

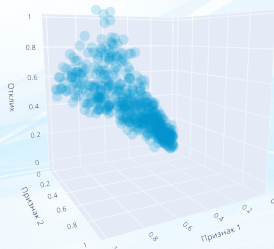


Признак 1

Нет, все хорошо!



Признак 2





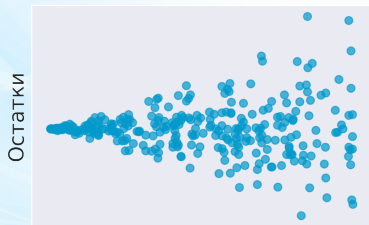
Критерии проверки на гомоскедастичность

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

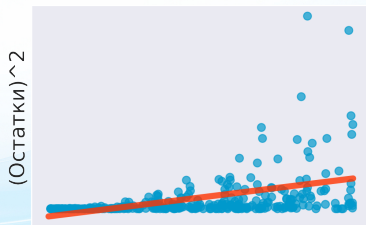
Критерий Бройша-Пагана

$R_{\hat{\varepsilon}^2}^2$ — коэф. детерминации для лин. регрессии предсказания $\hat{\varepsilon}^2$ по X

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



Признак



Признак



Критерии проверки на гомоскедастичность

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

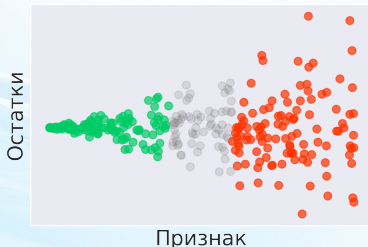
Критерий Голдфелда-Квандта

Упорядочим наблюдения по предполагаем. возрастанию дисперсий.

$\|Y - X\hat{\theta}_1\|^2$ — остатки регрессии по первым $\frac{n-r}{2}$ наблюдений, $r > 0$

$\|Y - X\hat{\theta}_2\|^2$ — остатки регрессии по последним $\frac{n-r}{2}$ наблюдений

$$\frac{\|Y - X\hat{\theta}_2\|^2}{\|Y - X\hat{\theta}_1\|^2} \sim F_{\frac{n-r}{2}-d, \frac{n-r}{2}-d} \quad \text{при } H_0$$





Что делать при гетероскедастичности?

- ▶ Если нужна только оценка θ — 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 λ подбирается по графику зависимости MSE от λ

- ▶ Использовать специальные оценки дисперсии, устойчивые к гетероскедастичности.

Устойчивые оценки дисперсии

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

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

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

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

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

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

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

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

3. Наличие автокорреляций — более сложный случай,
при котором зависимы элементы выборки.

Используются кластерное представление ковариационной матрицы или модели временных рядов.



Оценки Уайта

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

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

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.



Как ее применять?

Если автокорреляции отсутствуют, то выполнена асимптотическая нормальность оценки коэффициентов

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

НСЕ дает состоятельную оценку на матрицу B :

$$n\hat{\Sigma} \xrightarrow{P} B.$$

Данный факт позволяет строить асимптотические дов. интервалы для коэффициентов моделей и таргета, а также критерий Вальда для проверки линейных гипотез $H_0: T\theta = \tau$.



ВСЁ!