



Статистика

DS-поток

Лекция 11



6.1 Гауссовская линейная модель



Доска



Доверительные интервалы

Величина	Интервал
σ	$\left(\sqrt{RSS(\hat{\theta}) / \chi_{n-d, 1-\alpha/2}^2}, \sqrt{RSS(\hat{\theta}) / \chi_{n-d, \alpha/2}^2} \right)$
Дов. интервал для размера шума в отклике	
θ_j	$\left(\hat{\theta}_j \pm T_{n-d, 1-\alpha/2} \cdot \hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}} \right)$
Дов. интервал для коэффициента перед j -м признаком	
$x_0^T \theta$	$\left(x_0^T \hat{\theta} \pm T_{n-d, 1-\alpha/2} \cdot \hat{\sigma} \sqrt{x_0^T (X^T X)^{-1} x_0} \right)$
Дов. интервал для среднего отклика на объекте x_0	
$x_0^T \theta + \varepsilon$	$\left(x_0^T \hat{\theta} \pm T_{n-d, 1-\alpha/2} \cdot \hat{\sigma} \sqrt{1 + x_0^T (X^T X)^{-1} x_0} \right)$
Предск. интервал для наблюдаемого отклика на объекте x_0	



Значим ли признак x_j ?

Гипотеза о незначимости коэффициента θ_j

$H_0: \theta_j = 0$ vs. $H_1: \theta_j \{<, \neq, >\} 0$

Критерий Стьюдента

$$T_j^0(X, Y) = \frac{\hat{\theta}_j}{\hat{\sigma} \sqrt{(X^T X)^{-1}_{jj}}} \stackrel{H_0}{\sim} T_{n-d}$$

$T_j^0(X, Y)$ — Т-статистика критерия

Для $H_1: \theta_j \neq 0$ критерий $\{|T_j^0(X, y)| > T_{n-d, 1-\alpha/2}\}$

Если H_0 не отвергается, то можно считать, что θ_j отклоняется от нуля статистически незначимо. Возможно, признак j стоит убрать.



Значима ли группа признаков?

Пусть $X = \begin{pmatrix} X_1 & X_2 \\ n \times (d-k) & n \times k \end{pmatrix}$, $\theta = \begin{pmatrix} \theta_1^T & \theta_2^T \\ (d-k) \times 1 & k \times 1 \end{pmatrix}^T$

$H_0: \theta_2 = 0$ — гипотеза о незначимости второй группы

Критерий Фишера

$$\frac{(RSS(\hat{\theta}_1) - RSS(\hat{\theta})) / k}{RSS(\hat{\theta}) / (n - d)} \stackrel{H_0}{\sim} F_{k, n-d},$$

где $RSS(\hat{\theta}) = \|Y - X\hat{\theta}\|_2^2$, $RSS(\hat{\theta}_1) = \|Y - X_1\hat{\theta}_1\|_2^2$

Значима ли регрессия вообще?

Берем $k = d - 1$:

$$\frac{R^2 / (d - 1)}{(1 - R^2) / (n - d)} \stackrel{H_0}{\sim} F_{d-1, n-d}, \quad \text{где } R^2 = 1 - \frac{RSS(\hat{\theta})}{\|Y - \bar{Y}\|_2^2}$$



Таблица в statsmodels

```
# Load data
In [4]: dat = sm.datasets.get_rdataset("Guerry", "HistData").data

# Fit regression model (using the natural log of one of the regressors)
In [5]: results = smf.ols('Lottery ~ Literacy + np.log(Pop1831)', data=dat).fit()

# Inspect the results
In [6]: print(results.summary())
```

OLS Regression Results

```
=====
Dep. Variable:          Lottery    R-squared:                0.348
Model:                  OLS        Adj. R-squared:           0.333
Method:                 Least Squares    F-statistic:             22.20
Date:                   Mon, 14 May 2018    Prob (F-statistic):       1.90e-08
Time:                   21:48:09    Log-Likelihood:          -379.82
No. Observations:         86    AIC:                     765.6
Df Residuals:             83    BIC:                     773.0
Df Model:                 2
Covariance Type:         nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
Intercept	246.4341	35.233	6.995	0.000	176.358	316.510
Literacy	-0.4889	0.128	-3.832	0.000	-0.743	-0.235
np.log(Pop1831)	-31.3114	5.977	-5.239	0.000	-43.199	-19.424



Разрыв мозга

OLS Regression Results

```
=====
Dep. Variable:                y      R-squared:                0.991
Model:                        OLS    Adj. R-squared:           0.991
Method:                        Least Squares    F-statistic:              2651.
Date:                          Mon, 11 Feb 2019    Prob (F-statistic):       2.68e-97
Time:                          03:50:00    Log-Likelihood:           101.97
No. Observations:              100    AIC:                      -195.9
Df Residuals:                   96    BIC:                      -185.5
Df Model:                       4
Covariance Type:               nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
x1	87.9447	93.569	0.940	0.350	-97.788	273.677
x2	-68.1838	88.958	-0.766	0.445	-244.764	108.396
x3	-86.4095	93.573	-0.923	0.358	-272.150	99.331
x4	68.3874	88.954	0.769	0.444	-108.184	244.959



Совместное применение крит. Фишера и Стьюдента

Возможны случаи разрыва мозга...

1. Фишер: группа признаков значима

Стьюдент: ни один признак из группы не значим

- ▶ отдельные признаки плохо объясняют y , но в совокупности ОК
один в поле не воин
- ▶ признаки в группе мультиколлинеарны
каждой твари по паре

2. Фишер: группа признаков не значима

Стьюдент: в группе есть значимые признаки

- ▶ незначимые признаки маскируют влияние значимых
иголка в стоге сена



6.2 Анализ остатков

Остатки

В качестве оценки ошибки 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

$$H_0: Ee_i = 0$$



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 ε_i рассмотрим остатки $e_i = Y_i - \hat{Y}_i$

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

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

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

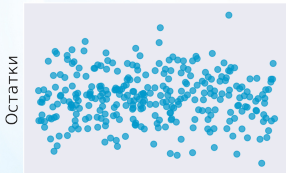
Проверять на однородность дисп. нужно **поправленные остатки**:

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



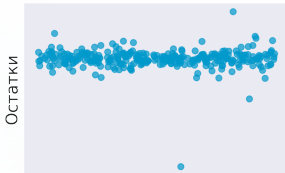
Визуальный анализ

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



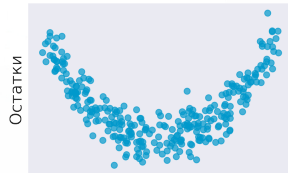
Признак

Все хорошо



Предсказание

Есть выбросы



Признак

Нужно добавить x^2



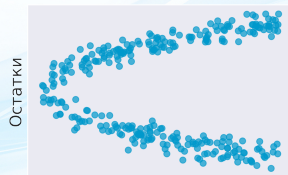
Признак

Гетероскедастичность



Номер наблюдения

Тренд



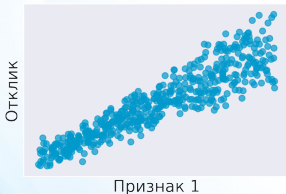
Признак

Неправильная модель

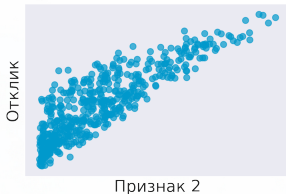


Визуальный анализ

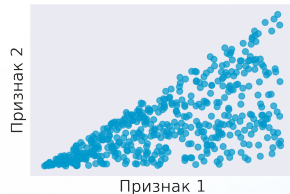
Что будет если строить графики зависимостей таргета от признаков:



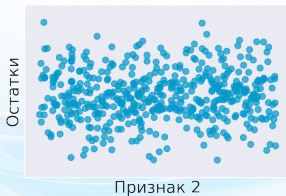
Гетероскедастичность?



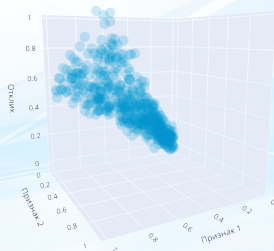
Гетероскедастичность?



Признаки зависимы



Нет, все хорошо!





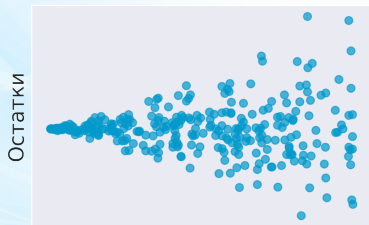
Критерии проверки на гомоскедастичность

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

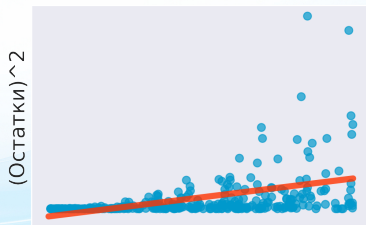
Критерий Бройша-Пагана

$R_{\hat{\varepsilon}^2}^2$ — 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

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

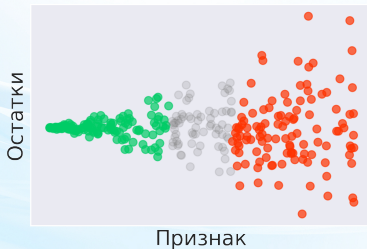
Критерий Голдфелда-Квандта

Упорядочим наблюдения по предполагаем. возрастанию дисперсий.

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 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 $RSS(\lambda)$ 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 — см. временные ряды.



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.



Асимптотическая нормальность при гетероскедаст.

Если автокорреляции отсутствуют, то

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

НСЕ дает состоятельную оценку на матрицу B :

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

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



ВСЁ!