

# ММРС | Регрессионный анализ

Если у нас не выполняются предпосыл. модели (прежде всего,  $E\varepsilon=0$ ,  $D\varepsilon=\sigma^2 I_n$ ,  $\varepsilon \sim N$ ,  $X^T X$  - обратима), то оценки коэфф. регрессии могут оказаться смещёнными и неустойч., доверит. интервалы будут иметь другие уровни доверия, т.е. всё нарушится. Поэтому надо уметь проверять эти предположения.

1)  $E\varepsilon=0$  - это к вводу адекватной модели относится, уже рассмотрели

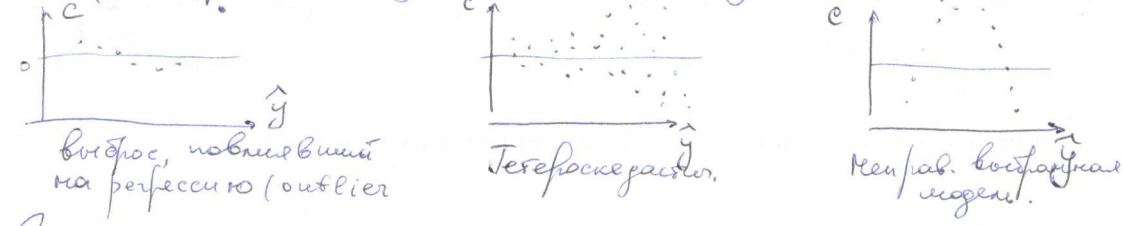
2)  $D\varepsilon=\sigma^2 I_n$ . Такая ситуация называется гомоскедастичностью, обратная, т.е. когда не все дисперсии равны, - гетероскедастичностью. Рассмотрим  $e_i = y_i - \hat{y}_i$ ,  $\hat{y}_i = (X(X^T X)^{-1} X^T y)_i$  - i-тая коэфф. оценки статистики. Если  $E\varepsilon=0$ , то  $E e_i = E(y - \hat{y}) = E(y - X\hat{\theta}) = X\hat{\theta} - X\hat{\theta} = 0$   
 $e = y - \hat{y} = (I - H)y$ , где  $H = X(X^T X)^{-1} X^T$  (эта часть матрицы называется "матрица остатков")

$Var e = Var(y - X\hat{\theta}) = Var((I - H)y) = (I - H) Var y (I - H)^T = \sigma^2 (I - H)(I - H)^T = \sigma^2 (I - H)$ , т.к.  $H^2 = H$  и  $H^T = H \Rightarrow D e_i = \sigma^2 (1 - H_{ii})$ .

Чтобы унифицировать все остатки, перейдём к нормир. остаткам:  $\frac{e_i}{\sqrt{D e_i}} = \frac{e_i}{\sigma \sqrt{1 - H_{ii}}}$   
 Но  $\sigma^2$  неизвестно  $\Rightarrow$  заменим  $\sigma^2$  на её оц.  $\frac{1}{n-k} \|y - X\hat{\theta}\|^2 = \frac{RSS}{n-k}$   
 получаем  $d_i = \frac{e_i}{\sqrt{\frac{RSS}{n-k}} \sqrt{1 - H_{ii}}}$  - стандартизированные остатки.

Замечание  $\sum_{i=1}^n h_{ii} = k$ , т.к.  $tr H = tr(X(X^T X)^{-1} X^T) = tr(X^T X(X^T X)^{-1}) = tr(I_k) = k$ , поэтому если  $k \ll n$ , то  $h_{ii}$  можно считать  $\approx 0$  и рассл.  $s_i = \frac{e_i}{\sqrt{\frac{RSS}{n-k}}}$  (стандартизир. остатки). Хотя стандартиз. остатки всё же зависимы, но при больших  $n$  их поведение похоже на незав.  $N(0,1)$ , (если предпол.  $\varepsilon \sim N(0, \sigma^2 I_n)$  верно).

Можно при анализе остатков пользоваться графическим анализом:



Так, проверка гетероскедастичности: проверим гипотезу  $H_0: \varepsilon \sim N(0, \sigma^2 I_n)$  (на самом деле, рассл.  $\varepsilon$  может отличаться от нормального)

## Критерий Таллфелда-Кванта

Упорядочим наблюд. по предполагаемому возрастанию дисперсий случай. ошибок, отбросим  $\alpha$  экстремальных наблюдений. Строим 2 регрессионные модели: по первым  $\frac{n-2}{2}$  наблюдениям и по  $\frac{n-2}{2}$  последним, вычислим по ним  $RSS_1$  и  $RSS_2$  - остат. суммы квадратов по каждой из моделей. Если  $H_0$  верна, то  $F = \frac{RSS_2}{RSS_1} \sim F_{\frac{n-2}{2}-k, \frac{n-2}{2}-k}$ , отвергаем, в противном случае  $F > \chi_{1-\alpha}$ , т.к.  $RSS_2 > RSS_1$ , по вводу модели



Критерий Зайта Рассл. вспомогат. модель  $e_i^2 = \alpha_0 + \sum_{j=1}^k \alpha_j x_{ij} + \sum_{j=1}^k \beta_j x_{ij}^2 + \nu_i$ , где  $e_i$  - обычные остатки модели и т.п.  $H_0: \alpha_j = \beta_j = 0 \forall j$ .

При верной  $H_0: nR^2 \sim \chi^2_{2k}$  (при большом кол-ве наблюд.), где  $R^2$  - коэфф. детерминации данной модели.

Есть ещё крит. Брайна-Папана, где  $\beta_j$  сразу = 0, и  $nR^2 \sim \chi^2_k$  в этой модели.

3)  $\varepsilon \sim N$ , проверка остатков на нормальность см. крит. проверки нормальности, чаще всего (из-за зависимости остатков) модн. удовлетв. хорошим видом QQ-plot и критерием Фарка-Бера. Если остатки не нормальны, то можно применить трансформ. Бокса-Кокса

Замечание: Удаление выбросов с помощью расст. Кука:  $D_i := \frac{e_i e_i^2}{RSS(k+1) (1 - H_{ii})^2}$  большие порога (напр., 1, 3D,  $\frac{4}{n}$ ), то удаление это наблюдение.

4) Независимость - проверяется корреляционным анализом. Если вдруг кажется, что остатки образуют временной ряд, то можно проверить на автокоррелированность (см. крит. Дарбина-Уотсона, а также критерии из авторегрессии), это будет далее в курсе

Замечание. Если от гетероскедаст. избавиться не удастся (напр., с помощью взвешенного МНК - т.е. одни наблюдения на что-то увелич., другие - уменьшаем (имеются в виду ~~отклонения~~)), то можно использовать для оценки значимости признаков критерии, основанные на устойчивой оц. дисперсии: White's heteroscedasticity-consistent estimator (HCE):  $\hat{\Omega} = n(X^T X)^{-1} (X^T \text{diag}(e_i^2) X) (X^T X)^{-1}$ .

Если вместо  $e_i^2$  поставить обычную оценку дисперсии, то она не получится. То получится стандарт. оценка дисперсии  $\hat{\sigma}^2$ .

Модификации: вместо  $e_i^2$  ставим  $e_i^2 \cdot \frac{n}{n-k}$ ,  $\frac{e_i^2}{1-H_{ii}}$ ,  $\frac{e_i^2}{(1-H_{ii})^2}$  (Маккинлонта Зайта) <sup>направки</sup>

5)  $X^T X$  - плохо обратима, это случается, когда столбцы  $X$  близки к линейной зависимости, т.е. сильная ~~корреляция~~ совместная корреляция между признаками. Такая ситуация называется мультиколлинеарностью.

а) В чём, собственно, проблема? Напомним, что  $\hat{\theta} \sim N(\theta, \sigma^2 (X^T X)^{-1})$ .

Если  $X^T X$  плохо обратима, то  $\det((X^T X)^{-1})$  очень большой, т.е. оценка  $\hat{\theta}$  будет крайне неустойчива. Критерием плохой обуслов.

матрицы  $X^T X$  авн. высокое знач. отношения  $\frac{\lambda_{\max}}{\lambda_{\min}}$  максимал. или минимал. собств. значений матрицы  $X^T X$ . Если  $\frac{\lambda_{\max}}{\lambda_{\min}}$  от 10 до 100, то уже стоит задуматься, а если  $> 1000$ , то ~~принимать меры~~.

б) Бывает так, что вектор-признак  $X_j$  имеет небольшое разброс значений и, как результат, будет коррелировать с  $X_0 = (1 \dots 1)^T$ .

$\Rightarrow$  следует перейти к центрированной матрице  $Z_j = \frac{x_{ij} - \bar{x}_j}{\sqrt{s_j^2}}$ , ну и центрировать отклик  $Y' = Y - X_0 \cdot \bar{y}$ .



Добавление к п. 4 - преобразование Бокса-Кокса.

Хотим, чтобы наши данные (манфистер, оталки) были похожи на нормальные. Пусть все инт. вектора  $Y$ -положит. Имеем  $Y \rightarrow Y - \min_i Y_i$

Делаем замену переменных:  $V = \begin{cases} \frac{Y^{\lambda-1}}{\lambda}, & \lambda \neq 0 \\ \ln Y, & \lambda = 0 \end{cases}$   $\bar{y} = (y_1, \dots, y_n)^{1/n}$  - среднее геометрич. отклика.

Находим опт.  $\lambda$  можно с помощью

а) критерий проверки нормальности

б) метода макс. пр-я

в) метода Бокса-Кокса: пусть  $\lambda \in (-a, a)$ , строим для каждого регрессию  $V$  по признакам  $X$ , получаем  $RSS(\lambda)$ , строим по нему график  $R^2(X)$ , выбираем таким образом оптимальн.  $\lambda$ . См. функцию `BoxCoxTrans()` пакета `car` в R.

~~Вывод, против адекватности~~ Возвращаемся к проблеме мультиколлинеар.

б) Удаление признаков с помощью VIF.

$$VIF(\theta_j) = \frac{1}{1 - R_j^2}, \text{ где } R_j^2 - \text{коэфф. детерминации модели}$$

$$X_j = c_0 + c_1 X_1 + \dots + c_{j-1} X_{j-1} + c_{j+1} X_{j+1} + \dots + c_k X_k.$$

Можно, что если  $R_j^2$  высок, то  $X_j$  объясняется <sup>хорошо</sup> ~~остальными~~ признаками, и его надо удалить. Обычно удаляют признак, если  $VIF_j > 4$ .

2) Регуляризаторы

• штрафная регрессия (ridge-regression):  $\hat{\theta}_{\text{ridge}} = \arg \min_{\theta} (\|Y - X\theta\|^2 + \lambda \|\theta\|^2)$ ,

$\hat{\theta}_{\text{ridge}} = (X^T X + \lambda I_n)^{-1} X^T Y$ , т.е. если  $\lambda > 0$ , то это будет обобщенн. о.н.к.

Проблема в том, что  $E\hat{\theta}_{\text{ridge}} = (I_k - \lambda W)\theta$ , где  $W = (X^T X + \lambda I_n)^{-1}$ , т.е. оценка смещенная.

Но зато  $\text{Var} \hat{\theta}_{\text{ridge}} = \sigma^2 W X^T X W$ , и с ростом  $\lambda$  она уменьшается, но смещение увеличивается.

Теорема  $\exists \lambda > 0$  такое, что  $E\|\hat{\theta}_{\text{ridge}} - \theta\|^2 < E\|\hat{\theta}_{\text{MNL}} - \theta\|^2$ , но все

на практике способа выбрать такое  $\lambda$  не придумали пока.

• лассо  $\hat{\theta}_{\text{LASSO}} = \arg \min_{\theta} (\|Y - X\theta\|^2 + \lambda \sum_{j=1}^k |\theta_j|)$ , и еще куда других регуляризаторов, см. Вайнерно или Mastie, Tibshirani

Регуляризаторы тесно связаны с байесовской постановкой задачи регрессии, где в кат. априорного распределения параметров модели  $\theta$  берется гауссовское (в ridge-regression) и распр. Лапласа (LASSO).

б) Ещё мы не обсудили проблему выбросов.

Выбросы очень сильно влияют на оценку н.к., см. квартиль Энскопта, напр.

Ранее для поиска оценки  $\hat{\theta}$  мы пользовались методом наим. квадратов.

$$\hat{\theta} = \arg \min_{\theta} \rho(Y - X\theta), \text{ где } \rho(\vec{v}) = \|\vec{v}\|$$

Так ведь можно взять другие  $\rho$ -уши  $\rho$ ! Требуем от них, чтобы они были 1) симметрич., 2) неотриц.; 3) монотонно  $\rho$ -неудовлетворяющим.

# Рефлексивный анализ | МСНС

(7)

(Обсуждаем робастные регрессионные модели). Пусть

Если  $\rho$  - дифф., то  $\psi = \rho'$  называется  $\psi$ -функцией влияния.

С пом. неё ~~фф-регрессии~~ ~~минимизация~~ поиск оценок  $\theta$  осуществл. из

такого ~~фф-мине~~ ~~фф-мине~~:  $\sum_{i=1}^n \psi(y_i - X_i^T \theta) \cdot X_i = 0$

Примеры  $\psi$ -функций  $\rho$ :

Тип	$L_2$	$L_1$	Huber	$L_p$	Cauchy	Грегман-Маклири
$\rho(x)$	$x^2/2$	$ x $	$\frac{x^2}{2} \cdot I( x  \leq k) + k( x  - \frac{k}{2}) \cdot I( x  > k)$	$ x ^p/p$	$\frac{c^2}{2} \log(1 + (\frac{x}{c})^2)$	$\frac{x^2/2}{1+x^2}$

Ну, требования: огранич.  $\psi$ -функция влияния (для устойчивости относ. выбросов) и единств. точка минимума.

$L_2$  -  $\psi$ -функция не ограничена,  $L_1$ , Коши, Маклуф - необязат. единств. решение.

Для Хьюбера обычно выбирают  $k = 1.345$ , кроме того, оценки методом

Хьюбера ~~похожи~~ ~~похожи~~ к асимпт. эффективными.

Ещё метод - метод скользящей медианы. Пусть  $m$  - фиксированное

число выбросов. Тогда заменяем  $y_i$  на  $\tilde{y}_i = \text{med}\{y_{i-m}, y_{i-m+1}, \dots, y_{i+m}\}$ .

Есть ещё метод, который, увы, плохо работает в ситуации, когда признаки зависимы - метод Тейбна:

есть регресс. модель  $y_i = \theta_0 + \sum_{j=1}^k \theta_j x_{ij} + \varepsilon_i$ , тогда берём

$\hat{\theta}_j = \text{med}\left\{\frac{y_i - y_k}{x_{ij} - x_{kj}}, i, k \in \{1, \dots, n\}\right\}$ ,  $\hat{\theta}_0 = \text{med}\left\{y_i - \sum_{j=1}^k \hat{\theta}_j x_{ij}\right\}$ . Зато устойчив к выбросам.