



Регрессия

**Регрессионный анализ по праву может
быть назван основным методом
современной математической статистики.**

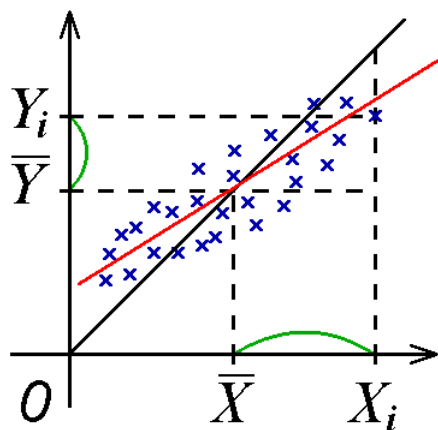
Н. Дрейнер, Г. Смит



Происхождение термина



Термин “регрессия” ввел Ф. Гальтон в своей статье “Регрессия к середине в наследовании роста” (1885 г.), в которой он сравнивал средний рост детей Y со средним ростом их родителей X (на основе данных о 928 взрослых детях и 205 их родителях). Гальтон заметил, что рост детей у высоких (низких) родителей обычно также выше (ниже) среднего роста популяции $\mu \approx \bar{X} \approx \bar{Y}$, но при этом отклонение от μ у детей меньше, чем у родителей. Другими словами, экстремумы в следующем поколении сглаживаются, происходит возвращение назад (*регрессия*) к середине.



По существу, Гальтон показал, что зависимость Y от X хорошо выражается уравнением

$$Y - \bar{Y} = (2/3)(X - \bar{X}).$$

«Живучесть» термина



В примечании переводчиков книги Дрейпер Н., Смит Г. *“Прикладной регрессионный анализ”* (кн. 1, с. 26) высказано интересное мнение по поводу “живучести” термина “регрессия”:

“Можно предположить, что его удивительная устойчивость связана с переосмыслением значения. Постепенно исходная антропометрическая задача, занимавшая Гальтона, была забыта, а интерпретация вытеснилась благодаря ассоциативной связи с понятием “регресс”, т. е. движение назад. Сначала берутся данные, а уж потом, задним числом, проводится их обработка. Такое понимание пришло на смену традиционной, еще средневековой, априорной модели, для которой данные были лишь инструментом подтверждения. Негативный оттенок, присущий понятию “регресс”, думается и вызывает психологический дискомфорт, поскольку воспринимается одновременно с понятиями, описывающими такой прогрессивный метод, как регрессионный анализ”.

Подгонка прямой

Пусть точки (x_i, η_i) получены в соответствии с моделью

$$\eta_i = a + bx_i + \varepsilon_i, \quad i = 1, \dots, n.$$

Здесь коэффициенты прямой a и b — неизвестные параметры, x_i — (неслучайные) значения переменной X , ε_i — независимые и одинаково распределенные случайные ошибки, $\mathbf{M}\varepsilon_i = 0$. Для нахождения оценок коэффициентов a и b применим метод наименьших квадратов (МНК).

Естественным условием точности подгонки *пробной* прямой $y = \alpha + \beta x$ служит близость к нулю всех *остатков*

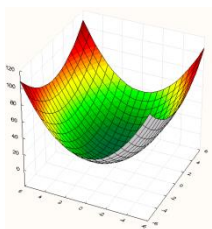
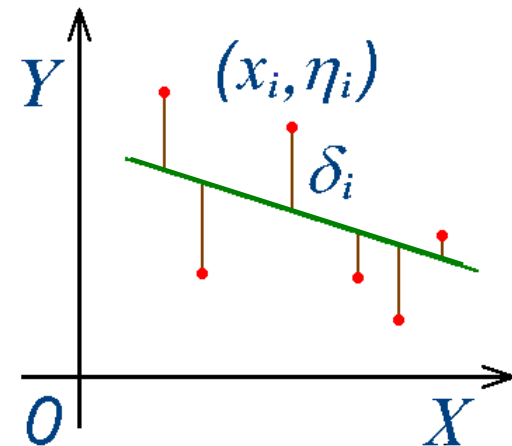
$$\delta_i(\alpha, \beta) = \eta_i - \alpha - \beta x_i.$$

Наиболее простые формулы для оценок \hat{a} и \hat{b} получаются, если в качестве *меры качества подгонки* взять

$$F(\alpha, \beta) = \sum_{i=1}^n \delta_i^2 = \sum_{i=1}^n (\eta_i - \alpha - \beta x_i)^2.$$

МНК-оценка (\hat{a}, \hat{b}) есть точка минимума функции $F(\alpha, \beta)$.

МНК был впервые опубликован Лежандром в 1805 г. Однако Гаусс утверждал, что он использовал МНК ещё в 1803 г.



Оценки коэффициентов a и b

Для вычисления МНК-оценок коэффициентов подгоняемой прямой используются следующие формулы:

$$\hat{b} = \sum_{i=1}^n (\eta_i - \bar{\eta})(x_i - \bar{x}) / \sum_{i=1}^n (x_i - \bar{x})^2, \quad \hat{a} = \bar{\eta} - \hat{b}\bar{x}.$$

Более устойчивым к выделяющимся наблюдениям является альтернативный метод оценивания коэффициентов a и b , предложенный Тейлом (Н. Theil) в 1950 г. Согласно этому методу оценки вычисляются по формулам

$$\begin{aligned} \tilde{b} &= MED \{ (\eta_j - \eta_i) / (x_j - x_i), \quad 1 \leq i < j \leq n \}, \\ \tilde{a} &= MED \{ \eta_i - \tilde{b}x_i, \quad i = 1, \dots, n \}. \end{aligned}$$

Причина устойчивости метода Тейла заключается в том, что одиночный «выброс» может исказить самое большее $(n - 1)$ оценку из $n(n - 1)/2$ оценок $(\eta_j - \eta_i) / (x_j - x_i)$ коэффициента наклона b .

Линейная регрессионная модель

Теперь рассмотрим зависимость признака η от $m \geq 2$ признаков X_1, X_2, \dots, X_m . Предположим, что эта зависимость (приблизительно) линейная в некотором диапазоне значений признаков с точностью до случайных ошибок. Иными словами, пусть

$$\eta = \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_m X_m + \varepsilon,$$

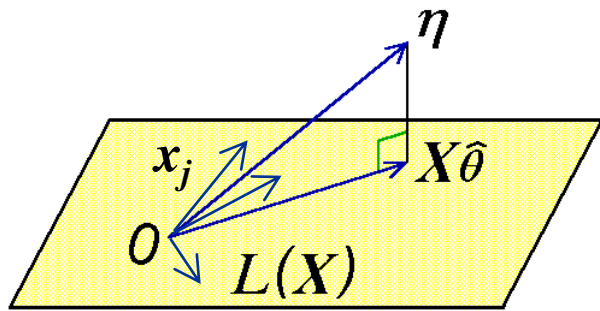
где $\theta_1, \theta_2, \dots, \theta_m$ — неизвестные коэффициенты (веса признаков), ε — случайная ошибка. Обозначим через x_{ij} значение признака X_j , где $j = 1, 2, \dots, m$ для объекта с номером i , где $i = 1, 2, \dots, n$. Тогда оценивание неизвестных коэффициентов методом наименьших квадратов заключается в минимизации по переменным $\theta_1, \theta_2, \dots, \theta_m$

$$F(\theta_1, \theta_2, \dots, \theta_m) = \sum_{i=1}^n (\eta_i - \theta_1 x_{i1} - \theta_2 x_{i2} - \dots - \theta_m x_{im})^2.$$

На следующем слайде объясняется, как эта задача сводится к простой процедуре — решению системы линейных уравнений.

Геометрическая интерпретация линейной регрессионной модели

Линейная модель



$$X\theta = \theta_1 x_1 + \dots + \theta_m x_m$$

$$\begin{pmatrix} \eta_1 \\ \eta_2 \\ \vdots \\ \eta_n \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nm} \end{pmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \\ \cdots \\ \theta_m \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\eta = X\theta + \varepsilon.$$

Вычисление МНК-оценок

$$X^T(\eta - X\hat{\theta}) = 0 \quad \text{или} \quad (X^T X) \hat{\theta} = X^T \eta.$$

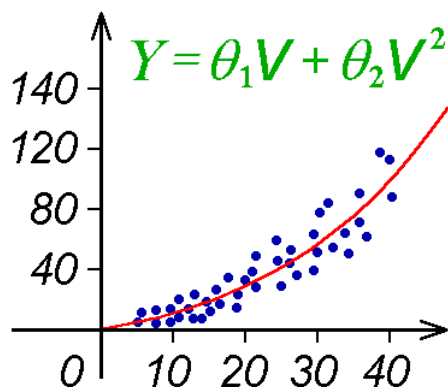
$$\hat{\theta} = (X^T X)^{-1} X^T \eta.$$

Система линейных уравнений относительно $\hat{\theta}$





Длина тормозного пути



На рисунке точками изображены результаты эксперимента по изучению зависимости между скоростью автомобиля V (в милях/час) и расстоянием Y (в футах), пройденным им после сигнала об остановке. Для каждого отдельного случая результат определяется в основном тремя факторами: скоростью V в момент подачи сигнала, временем реакции θ_1 водителя на этот сигнал и тормозами автомобиля. Автомобиль успеет проехать путь $\theta_1 V$ до момента включения водителем тормозов и еще $\theta_2 V^2$ после этого момента, поскольку согласно физическим законам теоретическое расстояние, пройденное до остановки с момента торможения, пропорционально квадрату скорости.

Таким образом, в качестве модели годится $Y = \theta_1 V + \theta_2 V^2$. Для экспериментальных данных были подсчитаны МНК-оценки $\hat{\theta}_1 = 0,76$ и $\hat{\theta}_2 = 0,056$. График параболы $Y = \hat{\theta}_1 V + \hat{\theta}_2 V^2$ приведен на рисунке.

Изучение бедности



В 30 американских округах были собраны следующие демографические характеристики:

POP_CHNG — Population change (1960-1970) [прирост населения];

N_EMPLD — No. of persons employed in agriculture [число жителей, занятых в сельском хозяйстве];

PT_POOR — Percent of families below poverty level [процент жителей, находящихся за чертой бедности];

TAX_RATE — Residential and farm property tax rate [местные налоги на землю и недвижимость];

PT_PHONE — Percent residence with telephones [доля телефонизации];

PT_RURAL — Percent rural population [доля сельского населения];

AGE — Median age [средний возраст жителей].

Изучим влияние на отклик **PT_POOR** остальных переменных (*предикторов*).

Визуальный анализ данных



- 1) Импортируйте файл `Poverty.txt` в `RStudio` под именем `p`
- 2) Постройте диаграммы размахов для всех признаков с помощью команды `boxplot`, а также для таблицы без 2-го столбца, и выявите явные «выбросы»
- 3) Замените явные «выбросы» на пропуски `NA` командой `ifelse`
- 4) Постройте матричную диаграмму рассеяния командой `plot` (нажмите `Zoom`, затем растяните окно во весь экран). Если обнаружите явные (т. е. не «хвостовые») двумерные «выбросы», то замените их на пропуски `NA`
- 5) Обратите внимание на 3-ю строку матричной диаграммы рассеяния и выясните, для каких предикторов наблюдается заметный наклон «облака» точек на диаграмме рассеяния предиктора с откликом `PT_POOR` (`p[,3]`)
- 6) Выясните, какие из 6 предикторов значимо на уровне 0,05 коррелируют с откликом `PT_POOR` (запишите результат функции `cor.test` в переменную `r` и узнайте `r$p.value`).
- 7) Выясните, какие предикторы значимо на уровне 0,05 (0,01) коррелируют между собой (напишите двойной цикл и запишите индикаторы `r$p.value<0.05` в матрицу `m`). Постройте на бумаге граф значимых связей, соединив рёбрами номера (или имена) значимо связанных предикторов

Решение

1) Импортируйте файл Poverty.txt в Rstudio под именем p

2) Выполните команды: `boxplot(p); boxplot(p[, -2])`

3) Можно удалить 1 или 2 «выброса», но не больше:

```
p[,2]=ifelse(p[,2]>10000, NA, p[,2]); p[,1]=ifelse(p[,1]>30, NA, p[,1])
```

4) На матричной диаграмме явных двумерных «выбросов» нет

5) Наклон «облака» наблюдается для признаков POP_CHNG, PT_PHONE, PT_RURAL (т. е. признаков с номерами 1, 5, 6)

6) Предикторы, указанные в пункте 5

7) Нажмите кнопку Install, установите пакет igraph, активируйте его, поставив «галку» перед ним на вкладке Packages. Выполните

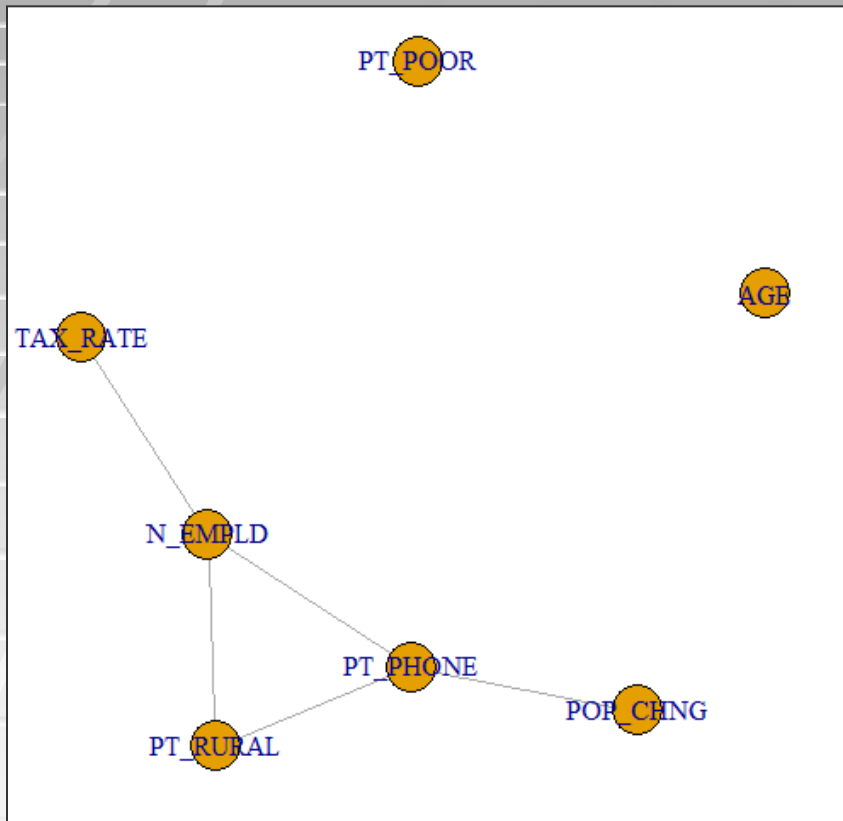
```
m=matrix(1:49, 7); s=names(p); rownames(m)=s; colnames(m)=s
for (i in 1:7) { for (j in 1:7) {
  r=cor.test(p[,i], p[,j])
  m[i, j]=ifelse(i!=j & i!=3 & j!=3 & r$p.value<0.05, 1, 0)
}}
```

затем постройте граф значимых связей между предикторами:

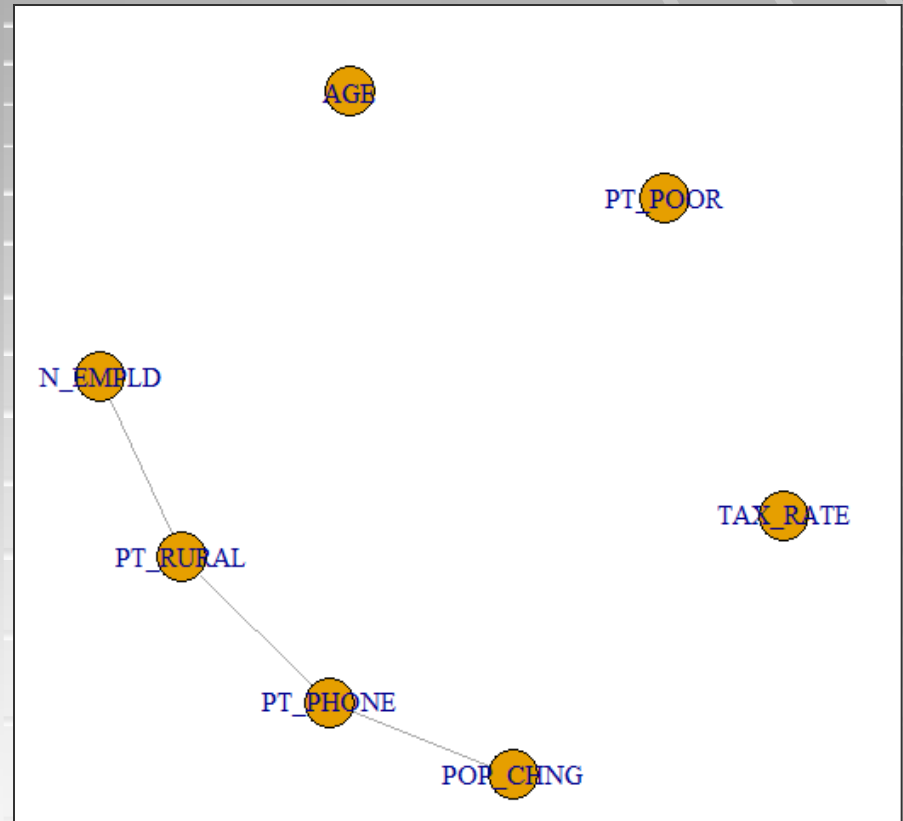
```
g=graph_from_adjacency_matrix(m, mode="undir"); plot(g) (см. ниже)
```

Графы значимых связей

P-value = 0,05



P-value = 0,01



Множественная регрессия

Для построения *линейной регрессионной модели* используйте функцию `summary(lm(p[,3]~., data=p[, -3]))`

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	30.772548	13.226947	2.327	0.03008
POP_CHNG	-0.323836	0.101604	-3.187	0.00443
N_EMPLD	0.002207	0.001078	2.047	0.05337
TAX_RATE	3.526140	3.537949	0.997	0.33027
PT_PHONE	-0.143406	0.136043	-1.054	0.30380
PT_RURAL	0.165268	0.062048	2.664	0.01454
AGE	-0.379517	0.260031	-1.460	0.15922

Estimate — коэффициенты линейной модели, *Intercept* — константа θ_0 , добавленная в модель, т. е. модель имеет вид

$$\eta = \theta_0 + \theta_1 X_1 + \dots + \theta_m X_m + \varepsilon$$

t value — характеристика степени влияния предиктора на отклик

Pr(>|t|) — фактические уровни значимости (отличия от нуля) соответствующих коэффициентов модели, которые вычисляются в предположении *нормальности распределения* ошибок наблюдений

Характеристики качества модели

Residual standard error: 3.382	on 21 degrees of freedom (2 observations deleted due to missingness)
Multiple R-squared: 0.7805	Adjusted R-squared: 0.7178
F-statistic: 12.45 on 6 and 21 DF,	p-value: 5.588e-06

Residual standard error — оценка для стандартного отклонения σ ошибки ε , которая вычисляется на основе регрессионных остатков δ_i по формуле

$$\sqrt{\sum \delta_i^2 / (n - m - 1)}$$

p-value — фактический уровень значимости всей модели, определяемый на основе F-статистики, который также вычисляется в предположении *нормальности распределения* ошибок наблюдений

Multiple R-squared — коэффициент детерминации (см. следующем слайд).

Коэффициент детерминации

Основным показателем силы связи между откликом и всеми предикторами служит *коэффициент детерминации*

$$R^2 = 1 - \frac{\sum_{i=1}^n (\eta_i - \tilde{\eta}_i)^2}{\sum_{i=1}^n (\eta_i - \bar{\eta})^2}, \quad \text{где } \tilde{\eta} = \mathbf{X}\hat{\theta}.$$

Таким образом, чем ближе R^2 к 1, тем в большей степени предикторы определяют (детерминируют) отклик. Если же значение R^2 близко к 0, то сглаживание наблюдений константой (их выборочным средним) мало отличается от сглаживания с помощью наилучшей линейной функции от предикторов.

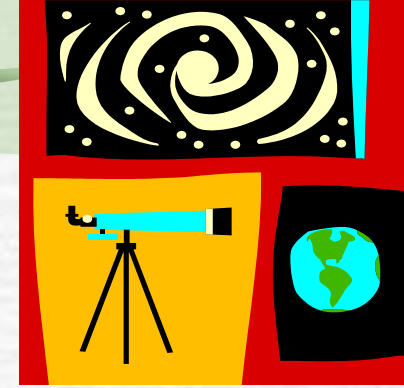
Показатель R^2 также называют *корреляционным отношением*, потому что он равен квадрату коэффициента корреляции R между откликом и *прогнозом*, который является проекцией отклика на пространство линейных комбинаций предикторов.

Избыточность (redundancy)

Чтобы решить, какие именно из тесно связанных между собой предикторов следует оставить в модели полезно учесть следующее:

- а) если предиктор связан с откликом причинно-следственной связью (при которой, изменяя предиктор, можно управлять откликом), то, как правило, следует оставить в модели именно его;
- б) если характер связи неизвестен, то лучше оставить предиктор, имеющий наибольший коэффициент корреляции с откликом;
- в) *степень избыточности* предиктора можно охарактеризовать его коэффициентом детерминации R^2 на основе остальных предикторов.

Анализ остатков

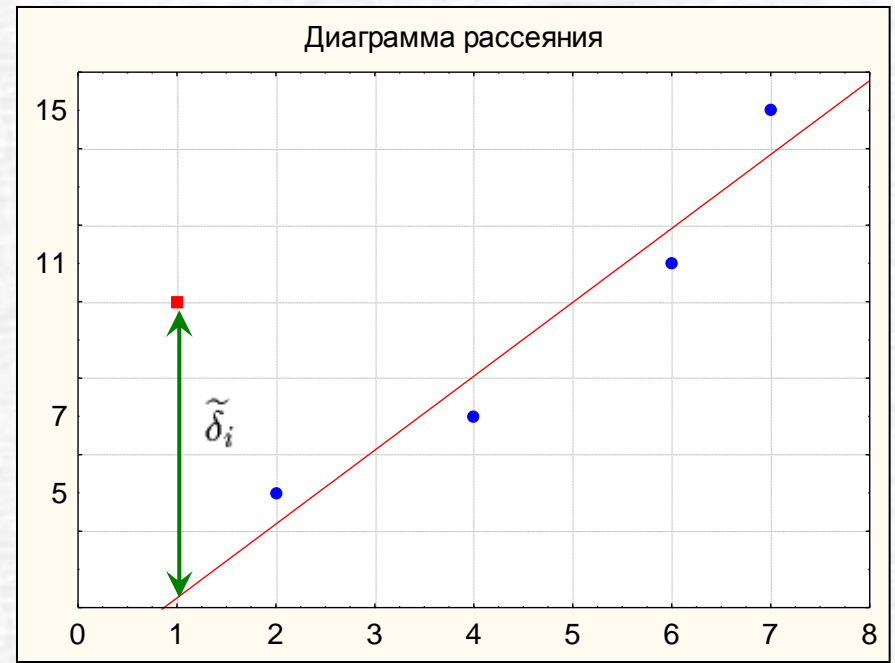
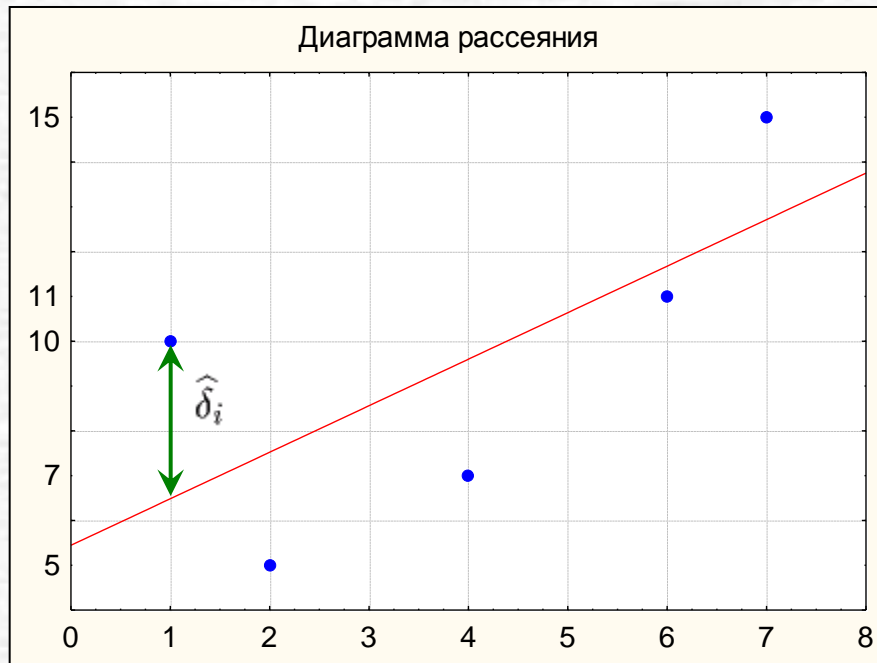


Почти все величайшие открытия в астрономии вытекают из рассмотрения того, что мы уже раньше назвали качественными или численными *остаточными феноменами*, иначе говоря, они вытекают из анализа той части числовых или качественных результатов наблюдения, которая «торчит» и остается необъясненной после выделения и учёта всего того, что согласуется со строгим применением известных методов.

Дж. Гершель, «Основы астрономии», 1849 г.

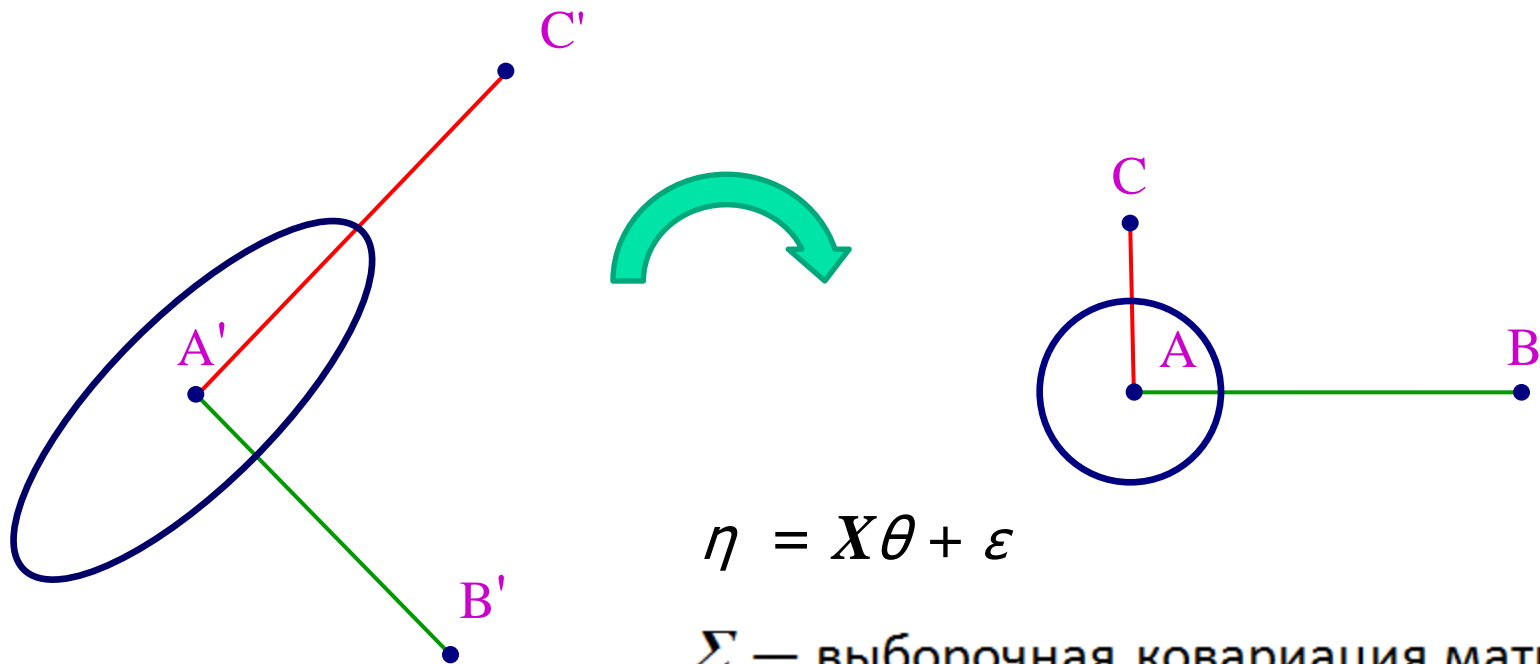
Deleted residuals

(остатки после удаления отдельных наблюдений)



Расстояние Махаланобиса

Предложено индийским статистиком Махаланобисом (Prasanta Chandra Mahalanobis) в 1936 году.

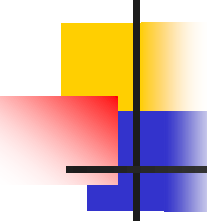


$$\eta = X\theta + \varepsilon$$

Σ — выборочная ковариация матрица столбцов матрицы X .

$$d^2(x, y) = (y - x)^T \Sigma^{-1} (y - x)$$

Расстояния Махаланобиса и Кука



Некоторые строки из таблицы данных могут оказаться «нетипичными по предикторам» в том смысле, что они не попадают внутрь 95%-доверительного эллипсоида рассеяния. Такие строки $\mathbf{x}_i = (x_{i1}, \dots, x_{im})$ будут иметь большие *расстояния Махаланобиса*

$$M_i = d^2(\mathbf{x}_i, \bar{\mathbf{x}}) = (\mathbf{x}_i - \bar{\mathbf{x}})\Sigma^{-1}(\mathbf{x}_i - \bar{\mathbf{x}})^T,$$

где $\bar{\mathbf{x}}$ — вектор-строка из выборочных средних всех предикторов, Σ — выборочная ковариация матрица предикторов, T — операция транспонирования. Эти строки, как рычаг (англ. *leverage*), способны отклонять регрессионную гиперплоскость от истинного положения.

Расстояние Кука (Cook's distance) определяются формулой

$$C_i = \left| \mathbf{X}\hat{\boldsymbol{\theta}}_{-i} - \mathbf{X}\hat{\boldsymbol{\theta}} \right|^2 / (m\hat{\sigma}^2),$$

где $\hat{\boldsymbol{\theta}}$ — вектор МНК-оценок, $\hat{\boldsymbol{\theta}}_{-i}$ — новый вектор МНК-оценок, получаемый при исключении i -й строки аналогично *deleted residuals*, $\hat{\sigma}$ — средняя ошибка. Расстояние Кука выражает степень влияния исключения i -й строки на изменение вектора прогноза $\mathbf{X}\hat{\boldsymbol{\theta}}$.

Критерий Дарбина — Уотсона

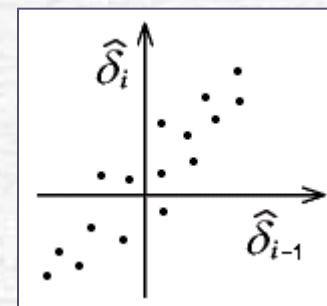
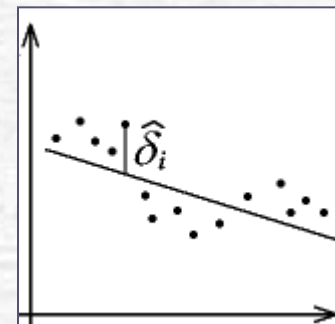
Данный критерий применяется для подтверждения значимости наблюдаемой сериальной корреляции остатков $\hat{\delta}_i$ (см. рисунок справа сверху). К подобному поведению остатков приводит справедливость следующей альтернативы H_1 для гипотезы H_0 независимости ошибок наблюдений ε_i :

$$H_1: \varepsilon_i = \rho \varepsilon_{i-1} + \zeta_i,$$

где $\zeta_i \sim \mathcal{N}(0, \sigma^2)$ и независимы, а $\rho \neq 0$, $|\rho| < 1$.

Статистикой критерия Дарбина — Уотсона служит

$$d = \sum_{i=2}^n (\hat{\delta}_i - \hat{\delta}_{i-1})^2 \bigg/ \sum_{i=1}^n \hat{\delta}_i^2.$$



H_0 отвергается

Неизвестно

H_0 принимается

0

d_L

d_U



Профессионализм

Я хотел бы спросить: «Что такое профессионал?» Многие, возможно, ответят, что профессионал — это человек, который очень много знает о своем предмете. Однако с этим определением я не мог бы согласиться, потому что никогда нельзя знать о каком-либо предмете действительно много. Я предпочёл бы такую формулировку: профессионал — это человек, которому известны грубейшие ошибки, обычно совершаемые в его профессии, и который, поэтому умеет их избегать.

В. Гейзенберг, «Физика и философия. Часть и целое».

«Ловушки» регрессии

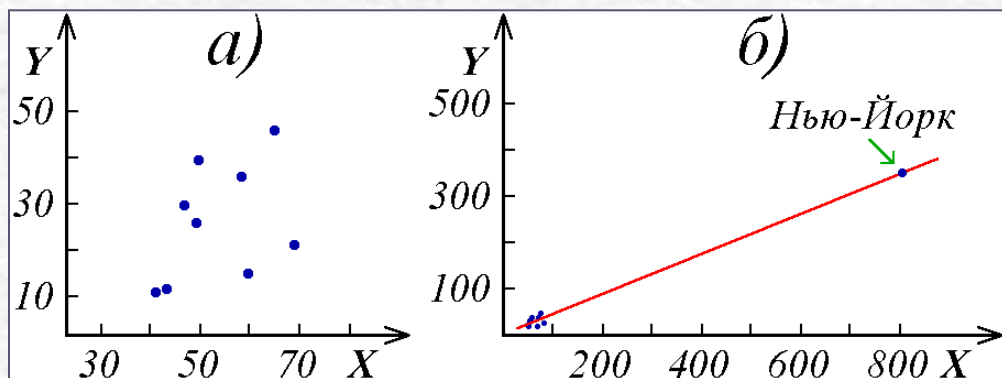


Существуют три вида лжи: ложь, наглая ложь и статистика.

Марк Твен

Есть несколько типичных ошибок (“тонких мест”), которые следует иметь в виду, применяя регрессионный анализ. Сами по себе, они достаточно очевидны. Тем не менее, о них часто забывают при работе с реальными данными и в результате приходят к неверным выводам.

- 1) Неоднородность данных
- 2) Коррелированность предикторов
- 3) Неадекватность модели
- 4) Скрытый фактор



«Ловушки» 2 – 4 подробно обсуждаются на следующих слайдах.



Плохая обусловленность матрицы

Ради краткости введём новые обозначения в уравнениях для вычисления МНК-оценок, полученных ранее:

$$B = X^T X, \quad d = X^T \eta.$$

Таким образом, поиск МНК-оценок сводится к решению системы линейных уравнений с матрицей коэффициентов B и правой частью d .

В случае **сильной коррелированности предикторов** матрица B оказывается *плохо обусловленной* (см. файл [Векторы и матрицы.pdf](#)). Для решений таких систем характерна катастрофическая неустойчивость к возмущению правой части. Классическим примером служит линейная система с *матрицей Гильберта* H , имеющей элементы $h_{ij} = 1 / (i + j - 1)$. Для численного эксперимента возьмём H размерности $m = 5$. Возмутим нулевую правую часть, положив последнюю компоненту вектора d равной **0,001**. В результате вместо нулевого решения получим

1	0,5	0,333333	0,25	0,2
0,5	0,333333	0,25	0,2	0,166667
0,333333	0,25	0,2	0,166667	0,142857
0,25	0,2	0,166667	0,142857	0,125
0,2	0,166667	0,142857	0,125	0,111111

 \times

0,63
-12,6
56,7
-88,2
44,1

 $=$

0
0
0
0
0,001

H d

Пошаговая (stepwise) регрессия

Этот метод используется в случае **большого числа** предикторов и в случае **коррелированности** предикторов.

Суть процедуры заключается в **постепенном увеличении** числа предикторов в модели. Опишем один из используемых подходов.

Сначала определяется предиктор, имеющий наибольший коэффициент корреляции с откликом. Если при его добавлении в модель коэффициент при нём значимо (скажем, на уровне **5%**) отличается от **0**, то предиктор оставляется.

На очередном шаге среди предикторов, ещё не включенных в модель, определяется тот, который имеет наибольшую *частную корреляцию* с откликом при устранении влияния предикторов, уже присутствующих в модели. Если при его добавлении в модель коэффициент при нём значимо отличается от **0**, то он оставляется.

Затем производится «чистка»: из модели удаляются все ранее включенные в нее предикторы, которые после добавления нового предиктора стали незначимыми (скажем, на уровне **10%**).

Добавление предикторов прекращается, когда на очередном шаге коэффициент при новом предикторе окажется незначимым.

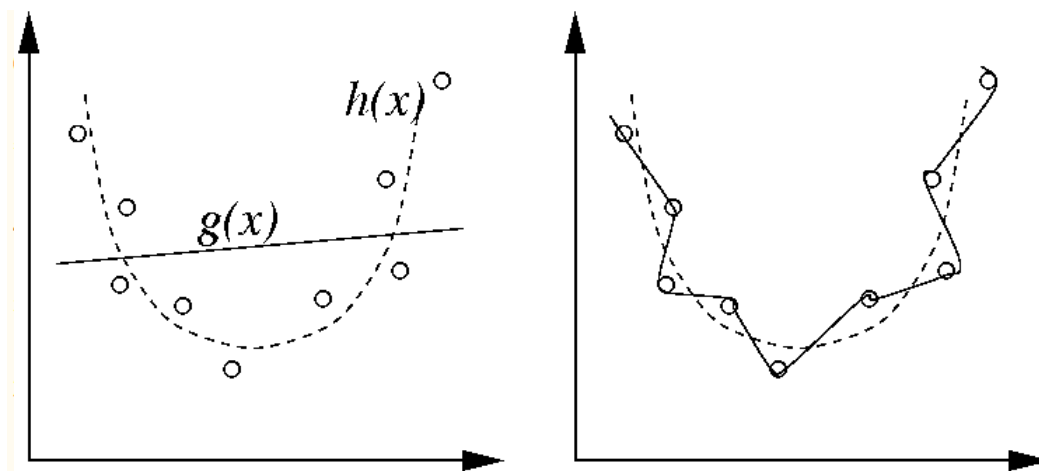
Альтернативный средством против коррелированности предикторов является переход к новым признакам — ортогональным главным компонентам.

Компромисс между смещением и дисперсией (bias-variance trade-off)

$$\mathbf{M}(\hat{\theta} - \theta)^2 = (\mathbf{M}\hat{\theta} - \theta)^2 + \mathbf{D}\hat{\theta}$$

Квадратичный риск = Квадрат смещения + Дисперсия

Регрессия



Штраф за сложность модели

Akaike's an Information Criterion (Akaike, 1973)

$$AIC = -2 \log L + 2K$$

Подробности см. в книге
Burnham K. P., Anderson R.
«Model Selection and
Multimodel Inference», 2002.

Здесь L — максим. правдоподобие, K — число параметров в модели.
Для линейной регрессионной модели с нормальными ошибками имеем

$$-2 \log L = n \log \hat{\sigma}^2, \quad \text{где} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \delta_i^2,$$

n — число наблюдений, δ_i^2 — регрессионные остатки, K включает в себя *intercept* (константу) и неизвестную дисперсию ошибок σ^2 .

Small sample (second order bias correction) version (Sugiura, 1978)

Для случая
 $n/K < 40$

$$AIC_c = -2 \log L + 2K + \frac{2K(K+1)}{n-K-1}$$

При $n > e^2 \approx 7,39$
штраф BIC
больше, чем
штраф AIC,
поэтому BIC-
модели проще

Bayesian Information Criterion (Schwarz, 1978)

$$BIC = -2 \log L + K \log n$$

Нелинейные преобразования



Поведение отклика	Уравнение	Усл. на b	x'	y'
Очень быстрый рост*)	$y = e^{a+bx}$	$b > 0$	x	$\ln y$
Быстрый (степенной) рост	$y = e^{a+b \ln x}$		$\ln x$	$\ln y$
Медленный рост				
Очень медленный рост				
Медленная стабилизация				
Быстрая стабилизация				
Кривая S-образной формы				

[Последняя функция на графике при $b > 0$ она возрастает, имея асимптоту $y = 1/a$ и перегиб в точке $x = -1/b$.
Переход к новым переменным сводит задачу к подгонке прямой $y = ax + b$ к данным $\ln y$ и $\ln x$.]

В книге Дж. Литлвуда «Математическая смесь» содержится любопытная классификация углов из книги по альпинизму: «Перпендикулярно — 60°, мой дорогой сэръ, абсолютно перпендикулярно — 65°, нависающе — 70°».

Содержательные модели

Модель	По всем наблюдениям		По части наблюдений	
	$\hat{\theta}$	$\hat{\sigma}$	$\hat{\theta}_{\text{тяж}}$	$\hat{\sigma}_{\text{лег}}$
1	$\hat{\theta}_1 = -984,7$ $\hat{\theta}_2 = 4,73$ $\hat{\theta}_3 = 4,70$	25,9	$\hat{\theta}_1 = 453,2$ $\hat{\theta}_2 = 0,62$ $\hat{\theta}_3 = -0,22$	81
2	$\hat{\theta}'_1 = 0,0011$ $\hat{\theta}'_2 = 1,556$ $\hat{\theta}'_3 = 1,018$	24,5	$\hat{\theta}'_1 = 266,4$ $\hat{\theta}'_2 = 0,203$ $\hat{\theta}'_3 = -0,072$	79
3	$\hat{\theta}_0 = 1,13 \cdot 10^{-4}$	26,6	$\hat{\theta}_0 = 1,11 \cdot 10^{-4}$	28



Влияние скрытого фактора

Скрытый фактор. Желание истолковывать регрессионную связь как причинно-следственную может приводить к парадоксам.

Во время второй мировой войны англичане исследовали зависимость *точности бомбометания* Z от ряда факторов, в число которых входили *высота бомбардировщика* H , *скорость ветра* V , *количество истребителей противника* X . Как и ожидалось, Z увеличивалась при уменьшении H и V . Однако (что поначалу представлялось необъяснимым), точность бомбометания Z возрастала также и при увеличении X .

Дальнейший анализ позволил понять причину этого парадокса. Дело оказалось в том, что первоначально в модель не был включен такой важный фактор, как Y — *облачность*. Он сильно влияет и на Z (уменьшая точность), и на X (бессмысленно высылать истребители, если ничего не видно). Сильные отрицательные причинно-следственные связи в парах (Y, Z) и (X, Y) привели к появлению положительного коэффициента при X в линейной регрессионной модели для Z .

Основные этапы регрессионного анализа

- 1) Выявление одномерных безусловных «выбросов» с помощью **диаграмм размахов** и их удаление [клавиша «Delete»]
- 2) Построение **диаграмм рассеяния** для поиска двумерных «выбросов», визуального изучения однородности данных и характера связи отклика с каждым предиктором в отдельности
- 3) Применение **монотонных преобразований** признаков в случае обнаружения нелинейной зависимости
- 4) Изучение **корреляционных связей** предикторов между собой и удаление избыточных предикторов из регрессионной модели
- 5) **Проверка значимости** всей линейной модели и каждого из коэффициентов модели в отдельности
- 6) Использование **пошаговой регрессии**, если имеет место коррелированность предикторов или их количество велико (во избежание перепогонки модели на обучающей выборке на один предиктор должно приходиться не менее 20 объектов)
- 7) **Анализ остатков**: скрытые «выбросы», адекватность модели
- 8) Выполнение **перепроверки** на случайной контрольной выборке

Домашнее задание

- 1) Импортируйте файл Poverty.txt в RStudio под именем `p`
- 2) Ради краткости введите переменные `n=nrow(p)`; `x=p[, -3]`; `y=p[, 3]`
- 3) Напишите программу для вычисления deleted residuals
(Чтобы набрать код программы, щёлкните по кнопке с белым крестиком внутри зелёного кружка, находящейся в левом верхнем углу экрана и выберите в меню пункт R Script Ctrl+Shift+N. В появившемся окне введите код программы. Для выполнения программы выделите весь код и нажмите кнопку Run с зелёной стрелкой, находящуюся над окном с кодом.)
- 4) Постройте линейную регрессионную модель `m` для самих `y` и `x`, затем постройте диаграмму рассеяния её регрессионных остатков и deleted residuals с номерами объектов (используйте команду `plot` с аргументами `asp=1` и `type="n"`, затем команду `text` с аргументом `labels=1:n`)
- 5) Выведите на диаграмму прямую с уравнением $y = x$ и нажмите кнопку Zoom. Объект из какой строки является явным «скрытым выбросом»?
- 6) Удалите из таблицы `p` строку этого объекта и повторите пункты 2-5. Есть ли другие объекты (строки таблицы), предположительно являющиеся «скрытыми выбросами»? Какие номера имеют соответствующие строки?