
Содержание

1	Методы получения оценок	2
1.1	Задачи	2
2	Свойства оценок	3
2.1	Практическое напоминание об условном ожидании и дисперсии	3
3	Асимптотические методы	5
4	Святая троица тестов	5
4.1	Чёрный трек	6
4.2	Тройка тестов в матричной форме	6
5	ШБ — МНК	7
5.1	Оптимизационная задача	7
5.2	Три простых модели	8
5.3	Как я перестал беспокоиться и полюбил матричное дифференцирование	11
5.4	Матричное представление регрессии	13
5.5	Решение оптимизационной задачи МНК с матрицами	14
5.6	Геометрия МНК	15
5.7	Показатели качества подгонки модели	18
5.8	Основные матрицы в линейной регрессии	20
5.9	Теорема Фриша — Во!	22
5.10	Кросс-валидация с выкидывание отдельных наблюдений	24
5.11	Задачи	25
6	Предпосылки о математическом ожидании и дисперсии	37
6.1	Иерархия зависимостей случайных величин	37
6.2	Ожидание и ковариационная матрица	38
6.3	Теорема Гаусса — Маркова для парной линейной регрессии	40
6.4	Теорема Гаусса — Маркова в общем виде	45
6.5	Статистические свойства остатков	48
6.6	Оценивание дисперсии	49
6.7	Неправильная спецификация модели	51
6.8	Задачи для семинара:	53
6.9	Компьютерные задачи для семинара:	61
6.10	Домашнее задание:	61
6.11	Чёрный трек:	61
7	Доверительные интервалы для коэффициентов	61
7.1	Случай многомерного нормального распределения	62
7.2	Независимость оценок β и $\hat{\sigma}^2$	64
7.3	Проверка гипотез о параметрах	64
8	Бутстрэп	64
9	Выбор функциональной формы	65

10 Гетероскедастичность	65
11 Мультиколлинеарность и метод главных компонент	65
12 Эндогенность	65
13 Эффекты воздействия	65
14 Задачи	65
15 Логистическая регрессия: точечные оценки	65
16 Логистическая регрессия: доверительные интервал	65
16.1 Смещение, цензурирование и ■■■■■■	66
16.2 Цензурирование	66
16.3 Усечение	66
16.4 Три осмысленных условных ожидания	67
Источники мудрости	68

1. Методы получения оценок

Методы получения оценок: метод максимального правдоподобия, метод моментов, метод наименьших квадратов.

1.1. Задачи

Задача 1. В оценках часто используют различные выражения с суммами. Пора перестать их бояться! Проверьте, верны ли следующие формулы,

а) $\sum_{i=1}^n (x_i - \bar{x}) = 0;$

б) $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x})x_i;$

в) $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i;$

г) $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i;$

д) $\sum_{i=1}^n x_i = n\bar{x};$

е) $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2;$

ж) $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y};$

з) $\sum_{i=1}^n x_i^2 = (\sum_{i=1}^n x_i)^2;$

и) $\sum_{i=1}^n x_i^2 = (n\bar{x})^2;$

к) $\sum_{i=1}^n \bar{x} = n\bar{x};$

л) $\sum_{i=1}^n x_i \bar{x} = n\bar{x}^2;$

$$\text{м) } \sum_{i=1}^n (x_i - \bar{x})y_i = 0.$$

Решение.

Здесь неотразимый ассист напишет решение

2. Свойства оценок

Свойства оценок: несмещённость, состоятельность, эффективность в классе.

2.1. Практическое напоминание об условном ожидании и дисперсии

Вспомним условное ожидание и условную дисперсию в дискретном случае:

Задача 2. Совместный закон распределения пары случайных величин (x, y) задан таблицей:

	$y = 1$	$y = 3$
$x = 1$	0.1	0.3
$x = 2$	0.1	0.1
$x = 4$	0.2	0.2

- а) Найдите $\mathbb{E}(y | x)$, $\text{Var}(y | x)$.
- б) Найдите $\mathbb{E}(y)$, $\mathbb{E}(x)$, $\text{Cov}(x, y)$, $\text{Var}(x)$.
- в) Найдите наилучшее линейное приближение $\text{BestLin}(y | x)$.

Решение. а) Условное математическое ожидание и дисперсия:

$$\mathbb{E}(y | x) = \begin{cases} 2.5, & x = 1, \\ 2, & x = 2, 4, \end{cases}$$

$$\text{Var}(y | x) = \begin{cases} 0.75, & x = 1, \\ 1, & x = 2, 4. \end{cases}$$

- б) Математические ожидания, ковариация и дисперсия:

$$\mathbb{E}(y) = 2.2, \quad \mathbb{E}(x) = 2.4, \quad \text{Cov}(x, y) = -0.28, \quad \text{Var}(x) = 1.84.$$

- в) Наилучшее линейное приближение:

$$\text{BestLin}(y | x) \approx 2.57 - 0.15 \cdot x.$$

Теперь вспомним, как считать условные характеристики случайных величин при наличии совместной плотности:

Задача 3. Пара случайных величин (x, y) имеет функцию плотности

$$f(x, y) = \begin{cases} (2x + 4y)/3, & \text{если } x \in [0, 1], y \in [0, 1], \\ 0, & \text{иначе.} \end{cases}$$

- а) Найдите $\mathbb{E}(y \mid x)$, $\text{Var}(y \mid x)$.
- б) Найдите $\mathbb{E}(y)$, $\mathbb{E}(x)$, $\text{Cov}(x, y)$, $\text{Var}(x)$.
- в) Найдите наилучшее линейное приближение $\text{BestLin}(y \mid x)$.

Решение.

Здесь мудрый ассист напишет решение

Особо обратим внимание на случай двумерного нормального распределения:

Задача 4. Пара случайных величин (x, y) имеет совместное нормальное распределение

$$\begin{pmatrix} x \\ y \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 3 \\ 2 \end{pmatrix}, \begin{pmatrix} 10 & -2 \\ -2 & 20 \end{pmatrix} \right)$$

- а) Найдите $\mathbb{E}(y \mid x)$, $\text{Var}(y \mid x)$.
- б) Найдите $\mathbb{E}(y)$, $\mathbb{E}(x)$, $\text{Cov}(x, y)$, $\text{Var}(x)$.
- в) Найдите наилучшее линейное приближение $\text{BestLin}(y \mid x)$.

Решение.

Здесь храбрый ассист напишет решение

Обратите внимание: для совместного нормального распределения условное ожидание $\mathbb{E}(y \mid x)$ и наилучшее линейное приближение $\text{BestLin}(y \mid x)$ идеально совпадают. Условное ожидание и линейное приближение совпадут и в том случае, если величина x принимает всего два значения. Убедимся в этом с помощью простой задачи

Задача 5. На первом шаге Илон Маск случайным образом выбирает одно из двух значений случайной величины x , $\mathbb{P}(x = 1) = 0.4$, $\mathbb{P}(x = 2) = 0.6$. На втором шаге Шивон Зилис выбирает значение y из экспоненциального распределения x с интенсивностью x .

- а) Найдите $\mathbb{E}(y \mid x)$, $\text{Var}(y \mid x)$.
- б) Найдите $\mathbb{E}(y)$, $\mathbb{E}(x)$, $\text{Cov}(x, y)$, $\text{Var}(x)$.
- в) Найдите наилучшее линейное приближение $\text{BestLin}(y \mid x)$.

Решение.

Здесь неотразимый ассист напишет решение

Теперь найдём условное ожидание и условную дисперсию для совместного нормального распределения в общем виде.

Определение 2.1 (наилучшее линейное приближение). Наилучшее линейное приближение величины r с помощью величины s — это линейная функция от s ,

$$\text{BestLin}(r \mid s) = \beta_1 + \beta_2 s,$$

где константы β_1 и β_2 находятся из решения задачи оптимизации $\mathbb{E}((r - \text{BestLin}(r, s))^2) \rightarrow \min_{\beta_1, \beta_2}$. При решении задачи окажется

$$\beta_1 = \mathbb{E}(r) - \beta_2 \mathbb{E}(s), \quad \beta_2 = \frac{\text{Cov}(r, s)}{\text{Var}(s)}$$

Определение 2.2 (линейно-независимые случайные величины). Величины r и s называются линейно-независимыми, если $\text{BestLin}(r \mid s) = \mathbb{E}(r)$.

Линейная независимость является симметричным явлением, $\text{BestLin}(r \mid s) = \mathbb{E}(r)$, если и только если $\text{BestLin}(s \mid r) = \mathbb{E}(s)$.

Задача 6. Выразите константы β_1 и β_2 в формуле для наилучшего линейного приближения

$$\text{BestLin}(r \mid s) = \beta_1 + \beta_2 s,$$

исходя из характеристик случайных величин r и s .

Решение. Выпишем целевую функцию в виде суммы

$$\mathbb{E}((r - \text{BestLin}(r, s))^2) = \text{Var}(r - \beta_1 - \beta_2 s) + (\mathbb{E}(r - \beta_1 - \beta_2 s))^2$$

Заметим, что β_1 не влияет на первое слагаемое, так как дисперсия константы равна нулю. И при этом, выбрав $\beta_1 = \mathbb{E}(r - \beta_2 s) = \mathbb{E}(r) - \beta_2 \mathbb{E}(s)$ мы добьёмся того, что второе слагаемое будет равно нулю, своему наименьшему возможному значению.

Остаётся минимизировать с помощью β_2 первое слагаемое.

$$\text{Var}(r - \beta_2 s) = \text{Var}(r) + \beta_2^2 \text{Var}(s) - 2\beta_2 \text{Cov}(r, s) \rightarrow \min_{\beta_2}.$$

Перед нами квадратичная функция от β_2 , следовательно,

$$\beta_2 = \frac{\text{Cov}(r, s)}{\text{Var}(s)}.$$

Обратите внимание, эта формула — родная «теоретическая» сестра «выборочной» формулы для парной регрессии

$$\hat{\beta}_2 = \frac{S_{xy}}{S_{xx}}.$$

Аналогия между оценкой и истинным коэффициентом действует и для первого коэффициента,

$$\beta_1 = \mathbb{E}(r) - \beta_2 \mathbb{E}(s), \quad \hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}.$$

И, попутно, мы замечаем, что условие $\text{Cov}(r, s) = 0$ равносильно условию $\text{BestLin}(r \mid s) = \mathbb{E}(r)$ или условию $\text{BestLin}(s \mid r) = \mathbb{E}(s)$.

3. Асимптотические методы

Центральная предельная теорема. Лемма Слуцкого. Дельта-метод. Построение асимптотических доверительных интервалов.

4. Святая троица тестов

Три классических теста: LM, LR, Wald.

Чёрный трек: тесты в матричной форме для вектора параметров?

4.1. Черный трек

4.2. Тройка тестов в матричной форме

Рассмотрим применение тестов W (тест Вальда), LR (тест отношения правдоподобия) и LM (тест множителей Лагранжа) для тестирования гипотез о параметрах модели.

Пусть требуется протестировать систему ограничений относительно вектора неизвестных параметров

$$H_0 : \begin{cases} g_1(\theta) = 0 \\ g_2(\theta) = 0 \\ \dots \\ g_r(\theta) = 0 \end{cases}$$

где $g_i(\theta)$ — функция, которая задаёт i -е ограничение на вектор параметров θ , $i = 1, \dots, r$.

Введём следующие обозначения:

$$\begin{aligned} \frac{\partial g}{\partial \theta^T} &= \begin{pmatrix} \partial g_1 / \partial \theta^T \\ \partial g_2 / \partial \theta^T \\ \vdots \\ \partial g_r / \partial \theta^T \end{pmatrix} = \begin{pmatrix} \frac{\partial g_1}{\partial \theta_1} & \frac{\partial g_1}{\partial \theta_2} & \dots & \frac{\partial g_1}{\partial \theta_k} \\ \frac{\partial g_2}{\partial \theta_1} & \frac{\partial g_2}{\partial \theta_2} & \dots & \frac{\partial g_2}{\partial \theta_k} \\ \dots & \dots & \dots & \dots \\ \frac{\partial g_r}{\partial \theta_1} & \frac{\partial g_r}{\partial \theta_2} & \dots & \frac{\partial g_r}{\partial \theta_k} \end{pmatrix} \\ \frac{\partial g^T}{\partial \theta} &= \begin{pmatrix} \frac{\partial g_1^T}{\partial \theta} & \frac{\partial g_2^T}{\partial \theta} & \dots & \frac{\partial g_r^T}{\partial \theta} \end{pmatrix} = \begin{pmatrix} \frac{\partial g_1}{\partial \theta_1} & \frac{\partial g_1}{\partial \theta_2} & \dots & \frac{\partial g_1}{\partial \theta_k} \\ \frac{\partial g_2}{\partial \theta_1} & \frac{\partial g_2}{\partial \theta_2} & \dots & \frac{\partial g_2}{\partial \theta_k} \\ \dots & \dots & \dots & \dots \\ \frac{\partial g_r}{\partial \theta_1} & \frac{\partial g_r}{\partial \theta_2} & \dots & \frac{\partial g_r}{\partial \theta_k} \end{pmatrix}, \quad \frac{\partial \ell}{\partial \theta} = \begin{pmatrix} \frac{\partial \ell}{\partial \theta_1} \\ \frac{\partial \ell}{\partial \theta_2} \\ \vdots \\ \frac{\partial \ell}{\partial \theta_k} \end{pmatrix} \\ I(\theta) &= -E \left(\frac{\partial^2 \ell}{\partial \theta \partial \theta^T} \right) = -\mathbb{E} \begin{pmatrix} \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} & \dots & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_k} \\ \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_2} & \dots & \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_k} \\ \dots & \dots & \dots & \dots \\ \frac{\partial^2 \ell}{\partial \theta_k \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_k \partial \theta_2} & \dots & \frac{\partial^2 \ell}{\partial \theta_k \partial \theta_k} \end{pmatrix} \end{aligned}$$

— информационная матрица Фишера

$\Theta_{UR} := \Theta$ — множество допустимых значений вектора неизвестных параметров без учёта ограничений

$\Theta_R := \{\theta \in \Theta : g(\theta) = 0\}$ — множество допустимых значений вектора неизвестных параметров с учётом ограничений

$\hat{\theta}_{UR} \in \Theta_{UR}$ — точка максимума функции ℓ на множестве Θ_{UR}

$\hat{\theta}_R \in \Theta_R$ — точка максимума функции ℓ на множестве Θ_R

Тогда для тестирования гипотезы H_0 можно воспользоваться одной из следующих ниже статистик:

$LR := -2(\ell(\hat{\theta}_R) - \ell(\hat{\theta}_{UR})) \stackrel{as.}{\sim} \chi_r^2$ — статистика отношения правдоподобия

$W := g^T(\hat{\theta}_{UR}) \cdot \left[\frac{\partial g}{\partial \theta^T}(\hat{\theta}_{UR}) \cdot I^{-1}(\hat{\theta}_{UR}) \cdot \frac{\partial g^T}{\partial \theta}(\hat{\theta}_{UR}) \right]^{-1} g(\hat{\theta}_{UR}) \stackrel{as.}{\sim} \chi_r^2$ — статистика Вальда

$LM := \left[\frac{\partial \ell}{\partial \theta}(\hat{\theta}_R) \right]^T \cdot I^{-1}(\hat{\theta}_R) \cdot \left[\frac{\partial \ell}{\partial \theta}(\hat{\theta}_R) \right] \stackrel{as.}{\sim} \chi_r^2$ — статистика множителей Лагранжа

5. ШБ — МНК

В этой главе мы приступим к изучению регрессионного анализа. Представьте, что вы работаете риэлтором или собираетесь приобрести квартиру в Москве. В обоих случаях Вас может интересовать, от чего зависит стоимость жилья (к примеру, от площади квартиры, расположения в том или ином районе Москвы и т.д.). Регрессионный анализ позволяет не только прогнозировать стоимость нового жилья, но и определять, какие факторы влияют на цену и в какую сторону — увеличивают или уменьшают её.

Для того чтобы корректно построить модель регрессии, сперва нам необходимо поставить оптимизационную задачу метода наименьших квадратов (МНК), найти её решение и обнаружить нестатистические свойства оценок.

5.1. Оптимизационная задача

Моделью парной регрессии называется модель вида

$$y_i = \beta_1 + \beta_2 x_i + u_i,$$

где y_i — значение зависимой (или, иначе, объясняемой) переменной для i -го наблюдения, x_i — значение объясняющей переменной (иногда её называют фактором или регрессором, иногда признаком :) для i -го наблюдения, β_1 — свободный коэффициент (константа), β_2 — коэффициент при факторе x , u_i — значение случайной ошибки для i -го наблюдения.

МНК-оценки коэффициентов модели парной регрессии находятся из решения задачи минимизации

$$Q(\hat{\beta}_1, \hat{\beta}_2) = \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) \rightarrow \min_{\hat{\beta}_1, \hat{\beta}_2}.$$

FOC (first order condition, условие первого порядка):

$$\begin{cases} \frac{\partial Q(\hat{\beta}_1, \hat{\beta}_2)}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) = 0, \\ \frac{\partial Q(\hat{\beta}_1, \hat{\beta}_2)}{\partial \hat{\beta}_2} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) x_i = 0. \end{cases}$$

$$\begin{cases} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) = 0, \\ \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) x_i = 0. \end{cases}$$

Разрешим данную систему уравнений относительно неизвестных $\hat{\beta}_1$ и $\hat{\beta}_2$, получим МНК-оценки коэффициентов β_1 и β_2 :

$$\begin{cases} \hat{\beta}_2^{\text{ols}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{sCov}(x, y)}{\text{sVar}(x)}, \\ \hat{\beta}_1^{\text{ols}} = \bar{y} - \hat{\beta}_2^{\text{ols}} \bar{x}. \end{cases}$$

Далее в целях снижения занудства мы иногда будем опускать нижний индекс ols и писать $\hat{\beta}$ вместо $\hat{\beta}_{\text{ols}}$.

Напомним, что $\text{sCov}(x, y)$ и $\text{sVar}(x)$ — это выборочная ковариация и выборочная дисперсия, определённые нами ранее в главе 2,

$$\begin{aligned} \text{sCov}(x, y) &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}, \\ \text{sVar}(x) &= \text{sCov}(x, x) = \frac{\sum (x_i - \bar{x})^2}{n - 1} \end{aligned}$$

Проиллюстрируем всё на конкретном примере: пусть $(x_1, y_1) = (1, 1); (x_2, y_2) = (1, 2); (x_3, y_3) = (0, 3); (x_4, y_4) = (0, 4)$. Тогда $\bar{y} = 2.5, \bar{x} = 0.5$.

$$\begin{cases} \hat{\beta}_2 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{0.5 \times (-1.5) + 0.5 \times (-0.5) + (-0.5) \times 0.5 + (-0.5) \times 1.5}{0.5^2 + 0.5^2 + (-0.5)^2 + (-0.5)^2} = \frac{-2}{1} = -2, \\ \hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x} = 2.5 - (-2) \times 0.5 = 3.5. \end{cases}$$

5.2. Три простых модели

Рассмотрим три простых регрессионных модели.

В качестве тренировки мы предлагаем вам вывести формулы оценок коэффициентов в этих моделях.

1. В регрессии на константу

$$y_i = \beta_1 + u_i$$

МНК-оценка параметра β_1 определяется по формуле

$$\hat{\beta}_1 = \bar{y}.$$

2. В модели регрессии без константы

$$y_i = \beta_1 x_i + u_i$$

МНК-оценка параметра β_1 равна

$$\hat{\beta}_1 = \sum_{i=1}^n \frac{x_i y_i}{x_i^2}.$$

В идеальной ситуации после оценивания коэффициентов модели вы ещё можете объяснить окружающим, какие полезные выводы можно сделать из получившихся оценок. Проинтерпретируем коэффициент при x_i . Увеличим объясняющую переменную на единицу, при этом старый прогноз равен $\hat{y}_i^{\text{old}} = \hat{\beta}_1 x_i$, новый прогноз равен $\hat{y}_i^{\text{new}} = \hat{\beta}_1 (x_i + 1)$ и разница прогнозов равна $\Delta y = \hat{y}_i^{\text{new}} - \hat{y}_i^{\text{old}} = \hat{\beta}_1$. Получается, что рост x_i на единицу приводит к изменению прогноза зависимой переменной на $\hat{\beta}_1$ единиц. Тот же самый результат можно получить взятием производной: $\frac{\partial \hat{y}_i}{\partial x_i} = \hat{\beta}_1$. Наша модель — линейная и изменение любого x_i на единицу окажет один и тот же эффект на предсказанный y_i .

Вариант, если CAPM уже пройдена на финансах

В качестве иллюстрации рассмотрим модель CAPM (Capital Asset Pricing Model) в её самом простом виде. В теории обычно предполагается, что премия за риск для некоторой ценной бумаги линейно зависит от рыночной премии за риск:

$$\text{premium}_i = \beta_1 \times \text{market-premium}_i + u_i$$

где premium_i — премия за риск¹ для ценной бумаги за i -й период, market-premium_i — рыночная премия за риск², β_1 — мера систематического (рыночного) риска бумаги (портфеля).

Найдём МНК-оценку для ценных бумаг в пищевом секторе: $\hat{\beta}_1 = 0.78$. Получается, что рост премии за риск на рынке на 1 пп (процентный пункт) приводит к тому, что доходность рассматриваемой ценной бумаги растёт на 0.78 пп. На рисунке 1 обратите внимание на то, что линия регрессии без

¹разница между доходностью бумаги и безрисковой ставкой, например, по государственным облигациям

²разница между доходностью рыночного портфеля, например, S&P500, и той же безрисковой ставкой

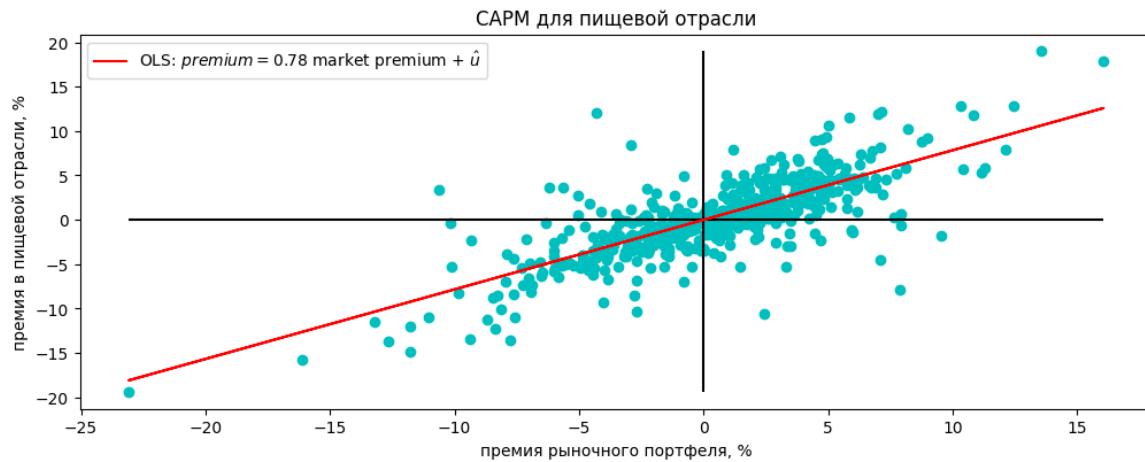


Рис. 1: CAPM модель, регрессия без константы

константы проходит через начало координат. Коэффициент получился меньше единицы, значит, покупать ценные бумаги в пищевом секторе — достаточно мало рисковая операция, так как рост премии за риск в пищевом секторе слабее, чем в среднем по рынку.

Вариант, если модель CAPM не пройдена

В качестве иллюстрации рассмотрим очень простую модель налогообложения, в которой будем для простоты игнорировать, что налоги могут считаться по прогрессивной шкале (так в целом было в России до 2025 года).

$$\text{taxes}_i = \beta_1 \times \text{income}_i + u_i$$

где taxes_i — объём подоходных налогов, уплачиваемый каждым индивидом, income_i — налогооблагаемая база, β_1 — ставка подоходного налога (для России до 2025 года ожидается, что при оценивании будет близка к 13%).

Казалось бы, такая линейная регрессия без константы должна подойти для нашей задачи, ведь если доходов нет, то и подоходные налоги платить не надо. Значит, линия регрессии должна проходить через начало координат.

3. В модели парной регрессии с константой

$$y_i = \beta_1 + \beta_2 x_i + u_i$$

МНК-оценки для параметров β_1 и β_2 выглядят следующим образом:

$$\hat{\beta}_2 = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{(x_i - \bar{x})^2}, \quad \hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}.$$

Результаты оценивания коэффициента $\hat{\beta}_2$ можно проинтерпретировать аналогично предыдущей модели, производная прогноза по регрессору равна $\frac{\partial \hat{y}_i}{\partial x_i} = \hat{\beta}_2$. После интерпретации эффекта, который x_i оказывает на y_i , можно заняться интерпретацией константы. Оценка коэффициента $\hat{\beta}_1$ — это прогноз \hat{y}_i при $x_i = 0$. Иногда оценка константы может иметь смысл. Давайте не поверим классической постановке модели и оценим модель CAPM с константой. В этом случае линия регрессии не обязана проходить через точку $(0, 0)$ на графике.

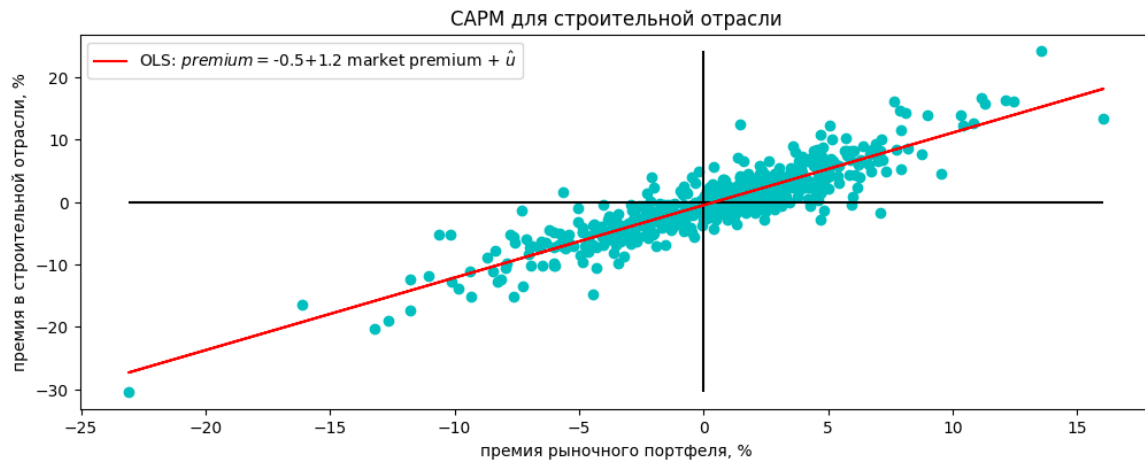


Рис. 2: CAPM модель, регрессия с константой

Будем использовать те же самые обозначения для премии за риск для выбранных ценных бумаг и рыночной премии за риск:

$$\text{premium}_i = \beta_1 + \beta_2 \times \text{market-premium}_i + u_i,$$

где β_2 – мера систематического (рыночного) риска бумаги (портфеля).

Оценённая линия регрессии для строительной отрасли получилась следующая:

$$\text{premium}_i = -0.5 + 1.2 \times \text{market-premium}_i + \hat{u}_i$$

Линия регрессии изображена на рисунке 2. При росте премии за риск для рыночного портфеля на 1 пп премия в строительной отрасли растёт на 1.2 пп, что говорит о, видимо, более высоких рисках в строительной сфере, которые хозяева фирм пытаются компенсировать. Более того, у нас оценена константа $\hat{\beta}_1 = -0.5$. Если бы доходность рыночного портфеля равнялась безрисковой ставке (рыночная премия за риск $\text{market-premium}_i = 0$), то премия за риск в строительной отрасли была бы -0.5 . То есть, если бы рыночный портфель был безрисковым, то инвесторы уходили бы из строительного сектора.

На самом деле, константа не всегда осмысленна. Рассмотрим короткий пример зависимости длины остановочного пути³ от скорости автомобиля:

$$\text{dist}_i = \beta_1 + \beta_2 \times \text{speed}_i + u_i,$$

Оценённая модель (изображена на рисунке 3) получилась следующей:

$$\text{speed}_i = -58.6 + 21.1 \times \text{dist}_i + \hat{u}_i.$$

Коэффициенты показывают, что изменение скорости на один километр в час (*ceteris paribus*, при прочих равных факторах) приводит к росту остановочного пути на 21 метр. Если же скорость равна нулю, то не стоит говорить, что остановочный путь должен составить -60 метров. В данном случае константа просто принимает значение параметра для наилучшего прохождения прямой через точки.

³путь, пройденный за время реакции водителя и фактического торможения

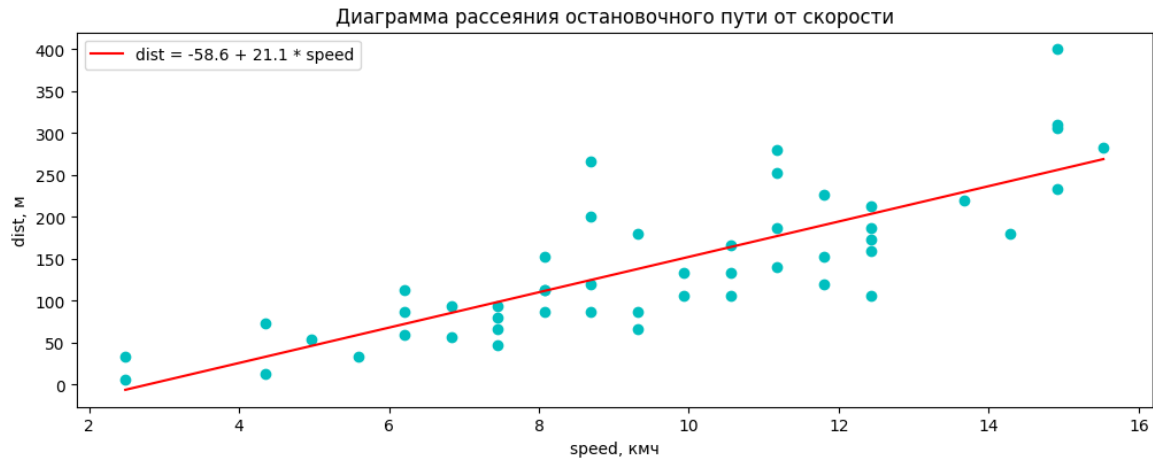


Рис. 3: Пример модели с константой, которая ничего не означает

5.3. Как я перестал беспокоиться и полюбил матричное дифференцирование

Мы планируем перейти к рассмотрению регрессий с большим числом факторов, поэтому полезна будет матричная запись модели и, соответственно, матричное дифференцирование. Этому и будет посвящён данный раздел.

Много полезных фактов про матрицы можно найти в шикарной книжке «Поваренная книга любителя матриц» Петерсона и Педерсона [PP12].

Нам чаще всего придётся дифференцировать скалярную функцию по векторному аргументу. По умолчанию вектор записывают столбцом и в большинстве источников производная по векторному аргументу тоже является столбцом,

$$\frac{\partial r}{\partial s} = \text{grad } r = \begin{pmatrix} \partial r / \partial s_1 \\ \partial r / \partial s_2 \\ \vdots \\ \partial r / \partial s_k \end{pmatrix}$$

Это сделано для того, чтобы размер результата дифференцирования совпадал с размером вектора s , по которому дифференцируют, и равнялся $[k \times 1]$. Например, для функции $f(x) = x_1^2 + x_2^3 + x_1 \cdot x_3^4$ векторная производная равна

$$\frac{\partial f}{\partial x} = \text{grad } f = \begin{pmatrix} 2x_1 + x_3^4 \\ 3x_2^2 \\ 4x_1x_3^3 \end{pmatrix}.$$

Сформулируем основные правила дифференцирования скалярных выражений по векторному аргументу:

$$\frac{\partial a^T s}{\partial s} = \frac{\partial s^T a}{\partial s} = \frac{\partial \sum s_i a_i}{\partial s} = a = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix}$$

$$\frac{\partial s^T A s}{\partial s} = \frac{\partial \sum_{ij} A_{ij} s_i s_j}{\partial s} = (A + A^T)s$$

В частности, для суммы квадратов правило превращается в

$$\frac{\partial s^T s}{\partial s} = \frac{\partial \sum s_i^2}{\partial s} = 2s = \begin{pmatrix} 2s_1 \\ 2s_2 \\ \vdots \\ 2s_n \end{pmatrix}$$

В более сложном случае дифференцирования вектора по вектору оказывается полезной матрица Якоби. В ней строки отвечают за элементы дифференцируемой функции, а столбцы — за элементы вектора, по которому дифференцируют.

Определение 5.1 (матрица Якоби). Для векторов r размера $[n \times 1]$ и s размера $[k \times 1]$ производной $\partial r / \partial s$ или матрицей Якоби называют матрицу, в которой дифференцируемые элементы записывают по строкам, а элементы, по которым берут производную, — по столбцам:

$$J = \frac{\partial r}{\partial s} = \begin{pmatrix} \partial r_1 / \partial s_1 & \partial r_1 / \partial s_2 & \dots & \partial r_1 / \partial s_n \\ \partial r_2 / \partial s_1 & \partial r_2 / \partial s_2 & \dots & \partial r_2 / \partial s_n \\ \vdots & \vdots & \ddots & \vdots \\ \partial r_n / \partial s_1 & \partial r_n / \partial s_2 & \dots & \partial r_n / \partial s_n \end{pmatrix}.$$

Например, для функции $f(x) = \begin{pmatrix} x_1 + 3x_2 \\ x_1 \cdot x_2^2 \end{pmatrix}$ матрица Якоби равна

$$J = \frac{\partial f}{\partial x} = \begin{pmatrix} 1 & 3 \\ x_2^2 & 2x_1 x_2 \end{pmatrix}.$$

Будьте бдительны! Обозначение для матрицы Якоби часто используют ровно то же, что для производной скалярной функции по векторному аргументу, $\partial r / \partial s$. При этом для скалярной функции r матрица Якоби «положит» все производные в строчку, а производная скалярной функции по векторному аргументу положит те же производные в столбик.

$$J = (\partial r / \partial s_1, \dots, \partial r / \partial s_n) \quad \frac{\partial r}{\partial s} = \begin{pmatrix} \partial r / \partial s_1 \\ \dots \\ \partial r / \partial s_n \end{pmatrix}$$

Избежать этой путаницы с упаковкой производных то в строчку, то в столбец позволяет дифференциал. У дифференциала всегда-всегда размерность совпадает с размерностью исходного дифференцируемого объекта.

Например, для функции $f(x) = \begin{pmatrix} x_1 + 3x_2 \\ x_1 \cdot x_2^2 \end{pmatrix}$ дифференциал равен

$$df = \begin{pmatrix} dx_1 + 3dx_2 \\ x_2^2 dx_1 + 2x_1 x_2 dx_2 \end{pmatrix}.$$

Если брать дифференциал скалярной функции s по векторному аргументу r , то он примет вид

$$ds = (\partial s / \partial r)^T dr = \text{grad}^T s dr,$$

поэтому из дифференциала легко извлечь градиент.

Например, если в результате долгих вычислений оказалось, что $ds = 2r^T dr$, то мы легко можем извлечь градиент-столбец $\text{grad } s = 2r$.

Напишем и правила для работы с дифференциалом. Здесь A, B — постоянные матрицы; a, b — постоянные векторы; R, S — матрицы переменных; r, s — векторы переменных.

Дифференциал постоянной матрицы равен нулю:

$$dA = 0.$$

При взятии дифференциала произведения важно сохранять порядок матриц R и S :

$$d(RS) = dR \cdot S + R \cdot dS.$$

В частности,

$$d(ARB) = A \cdot dR \cdot B.$$

Для суммы квадратов правило превращается в

$$ds^T s = d(\sum s_i^2) = 2s^T ds = 2 \sum s_i ds_i.$$

Линейность сохраняется для следа матрицы

$$d \text{trace } R = \text{trace } dR.$$

5.4. Матричное представление регрессии

Пусть теперь в модель для y включены k регрессоров x_1, x_2, \dots, x_k . Если в модель регрессии включена константа, то мы считаем, что $x_{i1} = 1$ для всех $i = 1, \dots, n$. Модель вида

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i$$

назовём моделью множественной регрессии.

Вектор зависимой переменной y имеет размер $[n \times 1]$, матрица признаков $X = [n \times k]$, вектор параметров модели $\beta = [k \times 1]$, вектор случайной ошибки $u = [n \times 1]$:

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}, \quad X = \begin{pmatrix} x_{11} & \dots & x_{1k} \\ x_{21} & \dots & x_{2k} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nk} \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_k \end{pmatrix}, \quad u = \begin{pmatrix} u_1 \\ u_2 \\ \dots \\ u_n \end{pmatrix}.$$

Используя введённые обозначения модель множественной регрессии

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i$$

можно переписать в матричном виде

$$y = X\beta + u.$$

Для наглядности можно преобразовать упражнение для парной линейной регрессии с константой из раздела 5.1 в матричный вид. Если при наличии в модели константы $x_{i1} = 1$ для всех i и $(x_{12}, y_1) = (1, 1)$; $(x_{22}, y_2) = (1, 2)$; $(x_{32}, y_3) = (0, 3)$; $(x_{42}, y_4) = (0, 4)$:

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix}, \quad X = \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix}.$$

5.5. Решение оптимизационной задачи МНК с матрицами

Для матричную регрессионной модели

$$y = X\beta + u$$

оптимизационную задачу МНК можно переписать как

$$Q(\hat{\beta}) = (y - X\hat{\beta})^T (y - X\hat{\beta}) = \hat{u}^T \hat{u} \rightarrow \min_{\hat{\beta}}.$$

Найдём МНК-оценку вектора β , используя матричное дифференцирование.

Запишем необходимое условие для задачи минимизации:

$$\frac{\partial Q(\hat{\beta})}{\partial \hat{\beta}} = -2X^T y + 2X^T X \hat{\beta} = 0$$

$$2X^T X \hat{\beta} = 2X^T y$$

Если матрица $X^T X$ размера $[k \times k]$ обратима, то формула для оценок принимает вид

$$\hat{\beta}_{\text{ols}} = (X^T X)^{-1} X^T y.$$

Если матрица $X^T X$ не обратима, то оценки МНК будут неединственными. Подобная проблемная ситуация возникает, если среди столбцов матрицы X есть линейно-зависимые, и называется строгой мультиколлинеарностью.

Для проверки достаточных условий второго порядка найдём матрицу Гессе в точке оптимума

$$\frac{\partial^2 Q(\hat{\beta})}{\partial \hat{\beta} \partial \hat{\beta}^T} = 2X^T X.$$

Если матрица $X^T X$ обратима, то она положительно определена и $\hat{\beta}_{\text{ols}}$ — точка глобального минимума.

Для полноты изложения найдём условие первого порядка для оптимального $\hat{\beta}$ и матрицу Гессе альтернативным способом — через дифференциал.

Запишем дифференциал для функции $Q(\hat{\beta})$ и выделим внутри него градиент,

$$dQ(\hat{\beta}) = 2(X\hat{\beta} - y)^T X d\hat{\beta} = (\text{grad } Q(\hat{\beta}))^T d\hat{\beta}.$$

Приравняем градиент к нулю

$$\text{grad } Q(\hat{\beta}) = 2X^T (X\hat{\beta} - y) = 0,$$

и получим прежнее условие первого порядка:

$$2X^T X \hat{\beta} = 2X^T y$$

Матрицу Гессе можно выделить внутри второго дифференциала,

$$d^2 Q = d(2(X\hat{\beta} - y)^T X d\hat{\beta}) = d\hat{\beta}^T \cdot 2X^T X \cdot d\hat{\beta}.$$

Как и ранее, матрица Гессе равна $\frac{\partial^2 Q(\hat{\beta})}{\partial \hat{\beta} \partial \hat{\beta}^T} = 2X^T X$.

Полученную оценку $\hat{\beta} = (X^T X)^{-1} X^T y$ называют МНК-оценкой $\hat{\beta}_{\text{ols}}$. Далее в целях снижения занудства мы иногда будем опускать нижний индекс ols и писать $\hat{\beta}$ вместо $\hat{\beta}_{\text{ols}}$.

Для нашего примера с

$$y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix}, \quad X = \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix}$$

$$X^T X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 4 & 2 \\ 2 & 2 \end{pmatrix}, (X^T X)^{-1} = \begin{pmatrix} 0.5 & -0.5 \\ -0.5 & 1 \end{pmatrix}, X^T y = \begin{pmatrix} 10 \\ 3 \end{pmatrix}.$$

$$\hat{\beta} = \begin{pmatrix} 0.5 & -0.5 \\ -0.5 & 1 \end{pmatrix} \begin{pmatrix} 10 \\ 3 \end{pmatrix} = \begin{pmatrix} 3.5 \\ -2 \end{pmatrix}.$$

Убедитесь, что результаты совпадают с примером в разделе 5.1.

5.6. Геометрия МНК

Кучу интересных геометрических фактов можно найти в статье [GD18] «Как встретились Гаусс, Марков и Пифагор?»

Метод наименьших квадратов имеет шикарную геометрическую интерпретацию. Геометрия позволяет не только вывести оценки, но и легко увидеть некоторые их свойства.

Для удобства рассмотрим случай двух регрессоров с константой

$$\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 a_i + \hat{\beta}_3 b_i.$$

Обозначим вектор из сплошных единиц буквой s , $s = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$. Буква s не является стандартным обозна-

чением для вектора из единиц и намекает на строевой лес, «ship forest». Также отметим, что вектор s нельзя называть единичным, единичный вектор — это вектор из нулей, в котором есть ровно одна единица.

С помощью вектора из единиц мы можем записать вектор прогнозов \hat{y} в виде

$$\hat{y} = \hat{\beta}_1 s + \hat{\beta}_2 a + \hat{\beta}_3 b.$$

Вектор \hat{y} — это линейная комбинация векторов s , a и b , $\hat{y} \in \text{span}(s, a, b)$. Под $\text{span}(s, a, b)$ мы обозначаем линейную оболочку векторов. Для наглядности можно представлять себе конкретный пример,

$$y = \begin{pmatrix} 2 \\ 2 \\ 3 \\ 4 \end{pmatrix}, \quad X = \begin{pmatrix} | & | & | \\ s & a & b \\ | & | & | \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 0 & 3 \\ 1 & 0 & 4 \end{pmatrix}$$

Теперь вспомним целевую функцию метода наименьших квадратов

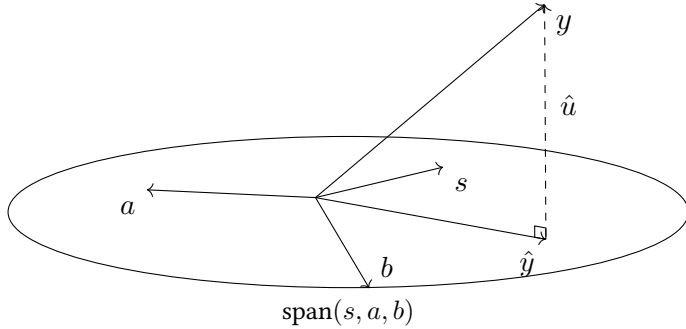
$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 \rightarrow \min_{\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3}.$$

Величина $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ — это квадрат длины вектора $y - \hat{y}$, то есть $\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \|y - \hat{y}\|^2$.

И теперь мы можем сформулировать метод наименьших квадратов геометрически!

Суть 5.1. Метод наименьших квадратов ищет вектор прогнозов \hat{y} внутри линейной оболочки регрессоров $\text{span}(s, a, b)$ поближе к вектору зависимой переменной y .

Конечно же, оптимальным решением оказывается проекция вектора y на линейную оболочку $\text{span}(s, a, b)$.



В случае произвольной матрицы X место линейной оболочки $\text{span}(s, a, b)$ займёт линейная оболочка всех столбцов матрицы X , $\text{span}(\text{col}_1 X, \text{col}_2 X, \dots, \text{col}_k X)$, которую мы для краткости будем обозначать $\text{colspan}(X)$.

Условием первого порядка будет ортогональность вектора остатков $\hat{u} = y - \hat{y}$ каждому из регрессоров,

$$\hat{u} \perp \text{span}(s, a, b) \quad \Leftrightarrow \quad \begin{cases} \hat{u} \perp s \\ \hat{u} \perp a \\ \hat{u} \perp b \end{cases}.$$

Поскольку оптимальный вектор прогнозов \hat{y} лежит в линейной оболочке $\text{span}(s, a, b)$, то вектор остатков \hat{u} перпендикулярен и ему тоже, $\hat{u} \perp \hat{y}$.

Условие ортогональности векторов означает, что скалярное произведение равно нулю, поэтому

$$\begin{pmatrix} -s^T \\ -a^T \\ -b^T \end{pmatrix} \cdot \hat{u} = 0 \quad \Leftrightarrow \quad X^T \hat{u} = 0.$$

А далее из условия ортогональности регрессоров и остатков $X^T \hat{u} = 0$ можно получить и явно формулы оценок всех коэффициентов. Подставим формулу для прогнозов, $\hat{y} = X\hat{\beta}$, и решим полученное уравнение $X^T(y - X\hat{\beta}) = 0$. Раскрываем скобки,

$$X^T y - X^T X \hat{\beta} = 0 \quad \Leftrightarrow \quad X^T X \hat{\beta} = X^T y.$$

Для рассматриваемого игрушечного примера,

$$X^T X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 2 & 3 & 4 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 0 & 3 \\ 1 & 0 & 4 \end{pmatrix} = \begin{pmatrix} 4 & 2 & 10 \\ 2 & 2 & 3 \\ 10 & 3 & 30 \end{pmatrix},$$

$$X^T y = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 2 & 3 & 4 \end{pmatrix} \begin{pmatrix} 2 \\ 2 \\ 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 11 \\ 4 \\ 31 \end{pmatrix},$$

Если нам повезло, и матрица $X^T X$ обратимая, то

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

$$\text{На нашем примере } \hat{\beta} = \begin{pmatrix} 4 & 2 & 10 \\ 2 & 2 & 3 \\ 10 & 3 & 30 \end{pmatrix}^{-1} \begin{pmatrix} 11 \\ 4 \\ 31 \end{pmatrix} = \begin{pmatrix} 12.75 & -7.5 & -3.5 \\ -7.5 & 5 & 2 \\ -3.5 & 2 & 1 \end{pmatrix} \begin{pmatrix} 11 \\ 4 \\ 31 \end{pmatrix} = \begin{pmatrix} 1.75 \\ -0.5 \\ 0.5 \end{pmatrix}$$

Заметим, что просто сократить матрицу X^T слева и справа нельзя потому, что она точно не обратимая! Матрица X^T имеет размер $[k \times n]$ и не является квадратной.

Готовая формула для вектора прогнозов равна

$$\hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T y.$$

$$\text{На нашем примере } \hat{y} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 0 & 3 \\ 1 & 0 & 4 \end{pmatrix} \begin{pmatrix} 1.75 \\ -0.5 \\ 0.5 \end{pmatrix} = \begin{pmatrix} 1.75 \\ 2.25 \\ 3.25 \\ 3.75 \end{pmatrix}.$$

Определение 5.2 (матрица-шляпница). Матрица $X(X^T X)^{-1} X^T$ часто обозначается буквой H и неформально называется матрицей-шляпницей (hat matrix). Она «надевает» на y шляпку, $Hy = \hat{y}$. Формально матрица H также называется матрицей-проектором. Она проецирует любой вектор на линейную оболочку всех регрессоров $\text{span}(\text{col}_1 X, \text{col}_2 X, \dots, \text{col}_k X)$.

$$\text{На нашем примере } H = \begin{pmatrix} 0.75 & 0.25 & 0.25 & -0.25 \\ 0.25 & 0.75 & -0.25 & 0.25 \\ 0.25 & -0.25 & 0.75 & 0.25 \\ -0.25 & 0.25 & 0.25 & 0.75 \end{pmatrix}. \text{ Заметим, что матрица-шляпница } H \text{ име-}$$

ет размер $[n \times n]$.

Вектор $\hat{\beta}$ состоит из k оценок $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$, поэтому его напрямую невозможно нарисовать в пространстве \mathbb{R}^n . Однако, оказывается, что веса, с которыми компоненты вектора зависимой переменной входят в каждую из оценок, можно изобразить! Заметим, что вектор оценок можно записать в виде

$$\hat{\beta} = (X^T X)^{-1} X^T y = W^T y, \text{ где } W = X(X^T X)^{-1}.$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_a \\ \hat{\beta}_b \end{pmatrix} = \begin{pmatrix} - & w_1^T & - \\ - & w_2^T & - \\ - & w_3^T & - \end{pmatrix} \cdot \begin{pmatrix} | \\ y \\ | \end{pmatrix} = \begin{pmatrix} | & | & | \\ w_1 & w_2 & w_3 \\ | & | & | \end{pmatrix}^T \cdot \begin{pmatrix} | \\ y \\ | \end{pmatrix} = \begin{pmatrix} 1.75 & -0.5 & -0.5 \\ -1.75 & 1.5 & 0.5 \\ 2.25 & -1.5 & -0.5 \\ -1.25 & 0.5 & 0.5 \end{pmatrix}^T \cdot \begin{pmatrix} | \\ y \\ | \end{pmatrix}.$$

То есть, каждая оценка $\hat{\beta}_j$ — это взвешенные наблюдения зависимой переменной y_1, y_2, \dots, y_n . Например, оценка первого коэффициента $\hat{\beta}_1$ — это скалярное произведение первого столбца W и зависимой переменной y ,

$$\hat{\beta}_1 = \langle \text{row}_1 W^T, y \rangle = \langle \text{col}_1 W, y \rangle = \langle w_1, y \rangle = w_{11}y_1 + w_{21}y_2 + \dots + w_{n1}y_n$$

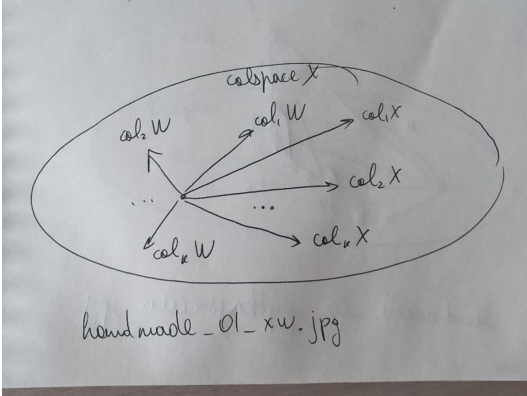
В нашем частном случае четырёх наблюдений,

$$\hat{\beta}_1 = 1.75y_1 - 1.75y_2 + 2.25y_3 - 1.25y_4.$$

Обратим внимание на запись $W = X(X^T X)^{-1}$. Для столбца весов w_1 она означает, что

$$\begin{pmatrix} | \\ w_1 \\ | \end{pmatrix} = \begin{pmatrix} 1.75 \\ -1.75 \\ 2.25 \\ -1.25 \end{pmatrix} = \begin{pmatrix} | & | & | \\ s & a & b \\ | & | & | \end{pmatrix} \cdot \text{col}_1(X^T X)^{-1} = \begin{pmatrix} | & | & | \\ s & a & b \\ | & | & | \end{pmatrix} \cdot \begin{pmatrix} 12.75 \\ -7.5 \\ -3.5 \end{pmatrix}$$

То есть, вектор весов $\text{col}_1 W$ лежит в линейной оболочке регрессоров, $\text{col}_1 W \in \text{colspan } X$. Например, при проецировании каждого столбца матрицы W на линейную оболочку регрессоров ничего не происходит, $HW = W$.



В этом равенстве можно убедиться и средствами линейной алгебры,

$$HW = X(X^T X)^{-1} X^T X (X^T X)^{-1} = X(X^T X)^{-1} = W.$$

5.7. Показатели качества подгонки модели

После оценивания регрессионной модели полезно проанализировать, насколько она «хороша». А именно, насколько похожи прогнозы \hat{y}_i на исходные наблюдения y_i . Для этого нужен показатель качества подгонки модели.

Назовём общей суммой квадратов (TSS) величину $\sum_{i=1}^n (y_i - \bar{y})^2$. Рассмотрим её разложение в сумму

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

Покажем, что $\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 0$:

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n (y_i - \hat{y}_i)\hat{y}_i - \bar{y} \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \hat{y}_i)\hat{y}_i - 0 = \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{\beta}_1 + \hat{\beta}_2 x_{i2} + \dots + \hat{\beta}_k x_{ik}) = \hat{\beta}_1 \sum_{i=1}^n (y_i - \hat{y}_i) + \hat{\beta}_2 \sum_{i=1}^n (y_i - \hat{y}_i)x_{i2} + \dots + \hat{\beta}_k \sum_{i=1}^n (y_i - \hat{y}_i)x_{ik} = 0, \end{aligned}$$

Здесь мы пользовались двумя фактами:

1. сумма остатков равна нулю $\sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n \hat{u}_i = 0$;
2. $\sum_{i=1}^n (y_i - \hat{y}_i)x_{ij} = 0$ для всех $j = 1, \dots, k$, так как остатки $\hat{u}_i = (y_i - \hat{y}_i)$ ортогональны регрессорам $x_j, j = 1, \dots, k$.

Заметим, что второе условие всегда выполнено, если используется метод наименьших квадратов. Первое условие будет гарантированно выполнено, если среди регрессоров будет присутствовать константа.

Таким образом, получаем разложение

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

Три составляющих этого разложения будут нам часто встречаться,

$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$	Общая сумма квадратов	Total Sum of Squares
$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	Объяснённая сумма квадратов	Explained sum of squares
$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$	Сумма квадратов остатков	Residual sum of squares

Давайте для нашего примера рассчитаем все три показателя. Имеем:

$$y = \begin{pmatrix} 2 \\ 2 \\ 3 \\ 4 \end{pmatrix}, \quad \hat{y} = \begin{pmatrix} 1.75 \\ 2.25 \\ 3.25 \\ 3.75 \end{pmatrix}, \quad \bar{y} = 2.75.$$

Тогда

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 = (2 - 2.75)^2 + (2 - 2.75)^2 + (3 - 2.75)^2 + (4 - 2.75)^2 = 2.75,$$

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = (1.75 - 2.75)^2 + (2.25 - 2.75)^2 + (3.25 - 2.75)^2 + (3.75 - 2.75)^2 = 2.5,$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (2 - 1.75)^2 + (2 - 2.25)^2 + (3 - 3.25)^2 + (4 - 3.75)^2 = 0.25.$$

Нетрудно убедиться, что выполняется тождество $TSS = ESS + RSS$.

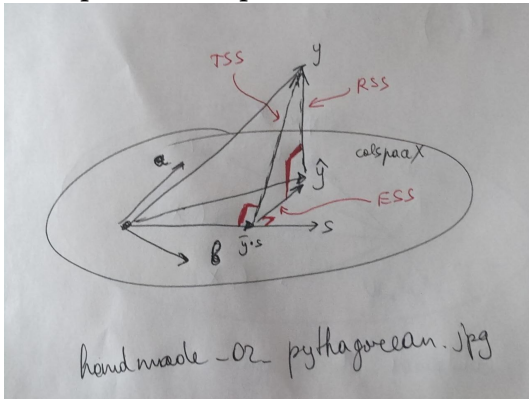
Будьте бдительны! Некоторые источники используют иные обозначения. По поводу TSS разногласий в литературе не возникает. Под RSS иногда понимают регрессионную сумму квадратов, regression sum of squares, то есть ESS в нашем курсе. Под ESS некоторые авторы подразумевают сумму квадратов остатков, error sum of squares, то есть наш RSS . Эта разница может быть критичной, например, при использовании формул из интернета, которые будут говорить противоположные нашему курсу вещи.

Запишем итоговую теорему.

Теорема 5.3. Если среди регрессоров присутствует константа, то

$$TSS = ESS + RSS.$$

На картинке это равенство окажется школьной теоремой Пифагора:



И в самом деле, $TSS = \sum (y_i - \bar{y})^2 = \|y - \bar{y} \cdot s\|^2$ — это квадрат длины гипотенузы $y - \bar{y} \cdot s$. Длина катета $y - \hat{y}$ в квадрате даёт $RSS = \sum (y_i - \hat{y}_i)^2 = \|y - \hat{y}\|^2$. Длина катета $\hat{y} - \bar{y} \cdot s$ в квадрате даёт $ESS = \sum (\hat{y}_i - \bar{y})^2 = \|\hat{y} - \bar{y} \cdot s\|^2$. Напомним, что s — это вектор строевого леса из единиц, а треугольник является прямоугольным в силу того, что мы проецируем вектор y на пространство $\text{colspan } X$.

Определение 5.4 (коэффициент детерминации). Коэффициентом детерминации называется статистика

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}.$$

При наличии константы в модели коэффициент детерминации $R^2 \in [0, 1]$ и показывает, какая доля разброса зависимой переменной относительно её среднего объясняется регрессионной моделью.

В нашем примере коэффициент детерминации равен

$$R^2 = \frac{ESS}{TSS} = \frac{2.5}{2.75} = 0.91.$$

Так как в нашем примере константа включена в число регрессоров, мы можем проинтерпретировать этот результат следующим образом: доля разброса зависимой переменной y относительно её среднего, объяснённая регрессионной моделью, составляет 0.91.

Ниже приведём альтернативное определение коэффициента детерминации.

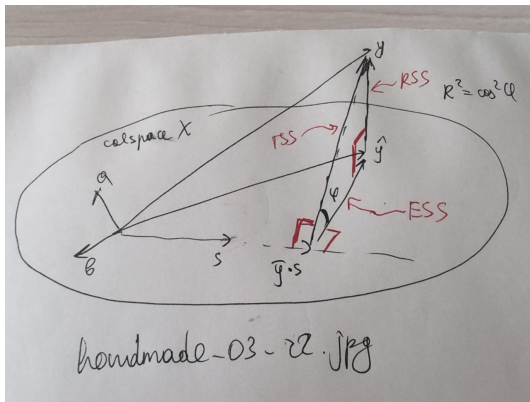
Определение 5.5 (коэффициент детерминации). В модели множественной регрессии с константой коэффициент детерминации R^2 — это квадрат выборочной корреляции между фактическими значениями зависимой переменной y и предсказанными значениями \hat{y} , полученными по модели,

$$R^2 = \text{sCorr}^2(y, \hat{y}).$$

При отсутствии константы в модели нарушается равенство $TSS = ESS + RSS$, в этом случае три формулы формулы $R^2 = ESS/TSS$, $R^2 = 1 - RSS/TSS$ и $R^2 = \text{sCorr}^2(y, \hat{y})$ дают разные результаты.

В регрессии без константы, в зависимости от используемой формулы, софт может выдать даже отрицательное значение R^2 . В разные пакеты могут быть встроены разные формулы для подсчёта R^2 без константы, например, в `sklearn` всегда используется формула $R^2 = 1 - RSS/TSS$, поэтому значения могут быть отрицательными. В `statsmodels` при отсутствии константы формула меняется: $R^2 = 1 - RSS / \sum_{i=1}^n y_i^2$ и значение всегда положительно.

На картинке R^2 окажется квадратом косинуса угла ϕ , $R^2 = \frac{ESS}{TSS} = \cos^2 \phi$.



5.8. Основные матрицы в линейной регрессии

Здесь мы приведём основные матрицы, используемые в линейной регрессии и их алгебраические свойства.

Вспомним определение матрицы-шляпницы:

Определение 5.6 (матрица-шляпница). Матрица $X(X^T X)^{-1} X^T$ часто обозначается буквой H и неформально называется матрицей-шляпницей (hat matrix). Она «надевает» на y шляпку, $H y = \hat{y}$. Формально матрица H также называется матрицей-проектором. Она проецирует любой вектор на линейную оболочку всех регрессоров $\text{span}(\text{col}_1 X, \text{col}_2 X, \dots, \text{col}_k X)$.

Теорема 5.7. Матрица-проектор $H = X(X^T X)^{-1} X^T$ обладает следующими свойствами:

- а) H — симметричная, $H^T = H$;
- б) H — идемпотентная, $H^2 = H$;
- в) $\text{rank } H = \text{trace } H = k$, где k — число столбцов матрицы X .

Доказательство. а) Симметричность, $H^T = H$:

$$H^T = X(X^T X)^{-1} X^T = H$$

- б) Идемпотентность, $H^2 = H$:

$$H^2 = X(X^T X)^{-1}(X^T X)(X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H$$

- в) $\text{rank } H = \text{trace } H = k$:

$$\text{trace } H = \text{trace}(X(X^T X)^{-1} X^T) = \text{trace}((X^T X)(X^T X)^{-1}) = \text{trace } I_k = k.$$

Здесь мы использовали свойство следа $\text{trace}(ABC) = \text{trace}(CAB)$.

□

Определим матрицу $M = I - H$. Матрица M , как и матрица H , является матрицей-проектором. Она проецирует любой вектор на ортогональное дополнение к линейной оболочке всех регрессоров $\text{span}(\text{col}_1 X, \text{col}_2 X, \dots, \text{col}_k X)$. Несложно убедиться, что матрица M так же, как и матрица H , симметричная и идемпотентная. При этом $\text{trace } M = \text{trace}(I_n - H) = \text{trace } I_n - \text{trace } H = n - k$. Проверьте симметричность и идемпотентность матрицы M самостоятельно.

Используя введенные матрицы, выразим вектор остатков в модели $y = X\beta + u$:

$$\hat{u} = y - \hat{y} = y - Hy = (I - H)y = My = M(X\beta + u) = Mu,$$

так как из геометрического смысла матрицы M следует, что $MX = 0$.

Пусть $s = (1 \ 1 \ \dots \ 1)^T$ — вектор размерности $[n \times 1]$, состоящий из единиц. Определим матрицу $\pi = s^T(s^T s)^{-1}s^T$. Матрица π — это матрица размерности $[n \times n]$ вида

$$\pi = \frac{1}{n} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \dots & \dots & \dots & \dots \\ 1 & 1 & \dots & 1 \end{pmatrix}.$$

В качестве домашнего упражнения покажите, что для произвольного вектора c размерности $[n \times 1]$ выполнено равенство $\pi \cdot c = \bar{c} \cdot s$, где \bar{c} — среднее арифметическое, посчитанное по элементам вектора c .

С помощью новых обозначений TSS , ESS и RSS могут быть записаны в матричном виде:

$$\begin{aligned} TSS &= (y - \bar{y})^T (y - \bar{y}) = (y - \pi y)^T (y - \pi y) = ((I - \pi)y)^T ((I - \pi)y) = y^T (I - \pi)^T (I - \pi)y = y^T (I - \pi)y, \\ ESS &= (\hat{y} - \bar{y})^T (\hat{y} - \bar{y}) = (Hy - \pi y)^T (Hy - \pi y) = (y(H - \pi))^T ((H - \pi)y) = y^T (H - \pi)^T (H - \pi)y = y^T (H - \pi)y, \\ RSS &= (y - \hat{y})^T (y - \hat{y}) = (y - Hy)^T (y - Hy) = ((I - H)y)^T ((I - H)y) = y^T (I - H)^T (I - H)y = y^T (I - H)y. \end{aligned}$$

5.9. Теорема Фриша — Во!

Сначала рассмотрим задачу. Джеймс Бонд для конспирации строит только регрессии на одну переменную. За один раз Джеймс может оценить ровно один коэффициент! Подобных регрессий он может построить сколь угодно много. Как может Джеймс Бонд, сохраняя конспирацию, оценить оба коэффициента в парной регрессии $\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i$?

Смотрите! На первом шаге Джеймс Бонд строит регрессию y на константу $s = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$. Он получает прогнозы первого шага $\hat{y}_i = \bar{y}$. Остатками первой регрессии являются $y_i^* = y_i - \bar{y}$.

На втором шаге Джеймс Бонд аналогично строит регрессию x на константу s и, аналогично, получает остатки $x_i^* = x_i - \bar{x}$.

И, наконец, на третьем шаге Джеймс строит регрессию полученных остатков y^* на остатки x^* , $\hat{y}_i^* = \hat{\beta} x_i^*$.

Оценка коэффициента на третьем шаге равна

$$\hat{\beta} = \frac{\sum y_i^* x_i^*}{\sum (x_i^*)^2} = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}.$$

Но именно это и есть оценка для $\hat{\beta}_2$ в парной регрессии $\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i$!

$$\hat{\beta}_2 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}.$$

Это совпадение не случайно! Оказывается множественную регрессию с любым количеством регрессоров можно разбить на несколько шагов со вспомогательными регрессиями с меньшим числом регрессоров.

Вместо непосредственного включения переменной x в качестве регрессора в модель можно сначала «очистить» от переменной x зависимую переменную y и остальные регрессоры, а затем оценить регрессию для «очищенных» переменных. Эту идею последовательных регрессий формализует теорема Фриша — Во — Ловелла.

Теорема 5.8 (Теорема Фриша — Во — Ловелла (англ. Frisch–Waugh–Lovell theorem, FWL theorem)). Рассмотрим два алгоритма.

Алгоритм А: оцениваем регрессию y на полный набор регрессоров X_1 и X_2 с помощью МНК:

$$\hat{y}_A = X_1 \hat{\beta}_1^A + X_2 \hat{\beta}_2^A.$$

Алгоритм В:

В1. Оцениваем регрессию y на часть регрессоров X_1 с помощью МНК:

$$\hat{y}_B = X_1 \hat{\beta}_1^B.$$

В2. Оцениваем регрессию каждого столбца из матрицы X_2 с помощью МНК:

$$\hat{X}_2^B = X_1 \hat{\beta}_1^B.$$

Уточним, что здесь $\hat{\beta}_1^B$ — это не вектор, а целая матрица, в которой содержатся оценки регрессии каждого столбца из матрицы X_2 на все регрессоры из матрицы X_1 .

В3. Определяем «очищенные» переменные как остатки регрессий первых двух шагов,

$$y^* = y - \hat{y}_B, \quad X_2^* = X_2 - \hat{X}_2^B.$$

В4. Оцениваем регрессию для «очищенных переменных»

$$\hat{y}^* = X_2^* \hat{\beta}_2^B.$$

Алгоритмы A и B дают одинаковые оценки коэффициентов $\hat{\beta}_2^A = \hat{\beta}_2^B$ и финальные векторы остатков $\hat{u}_A = y - \hat{y}_A = y^* - \hat{y}^* = \hat{u}_B$.

Доказательство. Определим матрицу-шляпницу H_1 , проецирующую на линейную оболочку столбцов блока X_1 , и матрицу $M_1 = I - H_1$, проецирующую на ортогональное дополнение к линейной оболочке столбцов X_1 ,

$$H_1 = X_1(X_1^T X_1)^{-1} X_1^T, \quad M_1 = I - H_1.$$

По определению, $H_1 X_1 = X_1$ и $M_1 X_1 = 0$.

Возьмём результат выполнения алгоритма A

$$y = X_1 \hat{\beta}_1^A + X_2 \hat{\beta}_2^A + \hat{u}_A$$

и домножим его на матрицу M_1 :

$$M_1 y = 0 \cdot \hat{\beta}_1^A + M_1 X_2 \hat{\beta}_2^A + M_1 \hat{u}_A$$

Заметим, что остатки \hat{u}_A алгоритма A ортогональны и регрессорам из блока X_1 , и регрессорам из блока X_2 . Сначала воспользуемся тем, что остатки \hat{u}_A уже лежат в подпространстве, ортогональном регрессорам блока X_1 . Дополнительное проецирование в это подпространство никак их не изменяет, $M_1 \hat{u}_A = \hat{u}_A$. Следовательно,

$$M_1 y = M_1 X_2 \hat{\beta}_2^A + \hat{u}_A$$

Теперь воспользуемся тем, что остатки \hat{u}_A уже лежат в подпространстве, ортогональном регрессорам блока X_2 , поэтому $X_2^T \hat{u}_A = 0$.

Ортогональны ли остатки \hat{u}_A и столбцы матрицы $M_1 X_2$? Проверим!

$$(M_1 X_2)^T \hat{u}_A = X_2^T M_1^T \hat{u}_A = X_2^T M_1 \hat{u}_A = X_2^T \hat{u}_A = 0.$$

Остаётся лишь сказать, что умножение на матрицу M_1 очищает переменные, $M_1 y = y^*$ и $M_1 X_2 = X_2^*$,

$$y^* = X_2^* \hat{\beta}_2^A + \hat{u}_A,$$

И мы видим идеальное совпадение с разложением алгоритма B ,

$$y^* = X_2^* \hat{\beta}_2^B + \hat{u}_B$$

В силу единственности разложения по ортогональному базису $\hat{\beta}_2^A = \hat{\beta}_2^B$ и $\hat{u}_A = \hat{u}_B$. □

5.10. Кросс-валидация с выкидывание отдельных наблюдений

Определение 5.9 (LOOCV). Кросс-валидация с поочередным выкидыванием отдельных наблюдений. На английском языке она часто сокращается LOOCV (leave one out cross validation).

Рассмотрим модель $y = X\beta + u$.

Оценим модель без первого наблюдения. Получим МНК-оценки $\hat{\beta}^{(-1)}$. С помощью этих оценок спрогнозируем первое наблюдение, получим прогноз \hat{y}_1^{CV} и ошибку прогноза \hat{u}_1^{CV} .

Вернём первое наблюдение в выборку и удалим второе наблюдение. Получим МНК-оценки $\hat{\beta}^{(-2)}$. С помощью этих оценок спрогнозируем второе наблюдение, получим прогноз \hat{y}_2^{CV} и ошибку прогноза \hat{u}_2^{CV} .

Поступим так с каждым наблюдением. На выходе получим вектор кросс-валидационных прогнозов \hat{y}^{CV} и вектор кросс-валидационных ошибок прогнозов $\hat{u}^{CV} = y - \hat{y}^{CV}$.

Теорема 5.10 (связь обычных и кросс-валидационных остатков). Если модель $y = X\beta + u$ оценивается с помощью МНК и проводится кросс-валидации с поочередным выкидыванием отдельных наблюдений, то:

$$\hat{u}_i = (1 - H_{ii}) \cdot \hat{u}_i^{CV},$$

где H — матрица-шляпница $H = X(X^T X)^{-1} X^T$, \hat{u} — остатки регрессии, а \hat{u}^{CV} — кросс-валидационные ошибки прогнозов.

Заметим, что сомножитель $(1 - H_{ii}) \in (0; 1)$. Другими словами, теорема численно формализует интуитивно ожидаемый результат: кросс-валидационные остатки по знаку совпадают с обычными остатками, а по абсолютной величине — больше, так как соответствующее наблюдение не используется при оценивании коэффициента.

Доказательство. Оценим модель без последнего наблюдения, $\hat{y}^d = X^d \hat{\beta}^d$. Буква d означает не степень, а удаление (deletion) последнего наблюдения, в частности, вектор y^d содержит $(n - 1)$ элемент, а матрица X^d имеет размер $[(n - 1) \times k]$.

Создадим вектор y^* , который будет отличаться от y только последним, n -м элементом: вместо настоящего y_n там будет стоять прогноз по модели без последнего наблюдения \hat{y}_n^d .

Раз уж мы добавили новую точку лежащую ровно на выборочной регрессии, то при оценки модели $\hat{y}^* = X\hat{\beta}^*$ мы получим в точности старые оценки $\hat{\beta}^* = \hat{\beta}^d$. Следовательно, и прогнозы эти две модели дают одинаковые, $\hat{y}_i^* = \hat{y}_i^d$.

А теперь посмотрим на последний элемент вектора $v = H(y^* - y)$.

С одной стороны, он равен последней строке матрицы H умножить на вектор $(y^* - y)$. В векторе $(y^* - y)$ только последний элемент ненулевой, поэтому $v_n = H_{nn}(\hat{y}_n^d - y_n)$.

С другой стороны, мы можем раскрыть скобки, и заметить, что $v = Hy^* - Hy$. И окажется, что $v_n = \hat{y}_n^* - \hat{y}_n = \hat{y}_n^d - \hat{y}_n$.

Отсюда

$$\hat{y}_n^d - \hat{y}_n = H_{nn}(\hat{y}_n^d - y_n)$$

Приводим подобные слагаемые и добавляем слева и справа y_n , получаем как раз то, что нужно:

$$y_n - \hat{y}_n = (1 - H_{nn})(y_n - \hat{y}_n^d)$$

□

5.11. Задачи

Задача 7. Оценим регрессию на константу $y_i = \beta_1 + u_i$ с помощью МНК. Найдите

- а) $\hat{\beta}_1$;
- б) \hat{y}_i ;
- в) ESS ;
- г) R^2 .

Решение. а) Задача минимизации для модели регрессии на константу:

$$Q(\hat{\beta}_1) = \sum_{i=1}^n (y_i - \hat{\beta}_1) \rightarrow \min_{\hat{\beta}_1}.$$

Условие первого порядка:

$$\frac{\partial Q}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_1) = 0.$$

Решаем уравнение:

$$\sum_{i=1}^n y_i - n\beta_1 = 0 \Rightarrow \hat{\beta}_1 = \frac{1}{n} \sum_{i=1}^n y_i.$$

б) Поскольку модель содержит только константу:

$$\hat{y}_i = \hat{\beta}_1 = \bar{y} \quad \text{для всех } i = 1, \dots, n.$$

в) Вычислим ESS :

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (\bar{y} - \bar{y})^2 = 0.$$

г) Вычислим компоненты:

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$ESS = 0$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = TSS$$

Коэффициент детерминации:

$$R^2 = \frac{ESS}{TSS} = 0.$$

Задача 8. Для модели парной линейной регрессии $y_i = \beta_1 + \beta_2 x_i + u_i$ получите оценки коэффициентов β_1, β_2 двумя путями:

а) используя готовую матричную формулу $\hat{\beta} = (X^T X)^{-1} X^T y$;

б) решая задачу оптимизации методом наименьших квадратов $\sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2 \rightarrow \min_{\hat{\beta}_1, \hat{\beta}_2}$.

Решение. а) Запишем модель в матричном виде:

$$y = X\beta + u, \quad \text{где} \quad X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

Вычислим необходимые матрицы:

$$X^T X = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}, \quad X^T y = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

Обратная матрица:

$$(X^T X)^{-1} = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix}$$

МНК-оценка:

$$\hat{\beta} = (X^T X)^{-1} X^T y = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i \\ n \sum x_i y_i - \sum x_i \sum y_i \end{pmatrix}$$

Таким образом:

$$\hat{\beta}_1 = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$
$$\hat{\beta}_2 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

б) Минимизируем сумму квадратов остатков:

$$Q(\hat{\beta}_1, \hat{\beta}_2) = \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2 \rightarrow \min_{\hat{\beta}_1, \hat{\beta}_2}.$$

Условия первого порядка:

$$\frac{\partial Q}{\partial \hat{\beta}_1} = -2 \sum (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) = 0,$$

$$\frac{\partial Q}{\partial \hat{\beta}_2} = -2 \sum x_i (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) = 0.$$

Получаем систему нормальных уравнений:

$$\begin{cases} n\hat{\beta}_1 + \hat{\beta}_2 \sum x_i = \sum y_i, \\ \hat{\beta}_1 \sum x_i + \hat{\beta}_2 \sum x_i^2 = \sum x_i y_i, \end{cases}$$

Решение системы:

$$\hat{\beta}_2 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2},$$
$$\hat{\beta}_1 = \frac{\sum y_i - \hat{\beta}_2 \sum x_i}{n} = \bar{y} - \hat{\beta}_2 \bar{x}.$$

Сравним результаты. Оба метода дают одинаковые выражения для оценок коэффициентов:

$$\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}$$
$$\hat{\beta}_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Задача 9. Покажите, что для модели парной линейной регрессии $y_i = \beta_1 + \beta_2 x_i + u_i$ выполняются следующие свойства:

- а) $\sum_{i=1}^n \hat{u}_i = 0$;
- б) $\sum_{i=1}^n \hat{y}_i = n\bar{y}$;
- в) $\sum_{i=1}^n \hat{u}_i x_i = 0$;
- г) $\bar{y} = \hat{\beta}_1 + \hat{\beta}_2 \bar{x}$

Решение. а) Оценки МНК получаются минимизацией суммы квадратов остатков:

$$Q(\hat{\beta}_1, \hat{\beta}_2) = \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2.$$

Условие первого порядка для $\hat{\beta}_1$:

$$\frac{\partial Q}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) = 0.$$

Отсюда:

$$\sum_{i=1}^n \hat{u}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = 0.$$

б) Из пункта а):

$$\sum_{i=1}^n \hat{u}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = 0 \implies \sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i.$$

Разделив на n , получим:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{\hat{y}}.$$

в) Условие первого порядка для $\hat{\beta}_2$:

$$\frac{\partial Q}{\partial \hat{\beta}_2} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) x_i = 0.$$

Отсюда:

$$\sum_{i=1}^n \hat{u}_i x_i = \sum_{i=1}^n (y_i - \hat{y}_i) x_i = 0.$$

г) Из пункта б):

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (\hat{\beta}_1 + \hat{\beta}_2 x_i) = \hat{\beta}_1 + \hat{\beta}_2 \bar{x}.$$

Задача 10. Для модели парной линейной регрессии $y_i = \beta_1 + \beta_2 x_i + u_i$

а) в скалярном виде выпишите TSS , RSS , ESS ,

б) покажите, что $TSS = RSS + ESS$.

Решение. а) Для модели парной линейной регрессии $y_i = \beta_1 + \beta_2 x_i + u_i$ найдём все компоненты:

- Общая сумма квадратов (Total Sum of Squares, TSS):

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2,$$

где $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ — выборочное среднее зависимой переменной.

- Объяснённая сумма квадратов (Explained Sum of Squares, ESS):

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2,$$

где $\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i$ — предсказанное значение y_i по оценённой модели.

- Остаточная сумма квадратов (Residual Sum of Squares, RSS):

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

где $y_i - \hat{y}_i = \hat{u}_i$ — остаток модели.

б) Докажем равенство $TSS = ESS + RSS$.

Начнём с определения TSS и разложим отклонение $y_i - \bar{y}$ на две составляющие:

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}).$$

Возведём обе части в квадрат и просуммируем по всем наблюдениям:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n [(y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})]^2.$$

Раскроем квадрат в правой части:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}).$$

Заметим, что:

- $\sum_{i=1}^n (y_i - \hat{y}_i)^2 = RSS$,
- $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = ESS$,
- перекрёстное произведение $\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 0$ (доказательство ниже).

Докажем, что $\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 0$. Заметим, что:

- $\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i$, где $\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}$ (из условий первого порядка МНК).
- $\sum_{i=1}^n \hat{u}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = 0$ (свойство остатков МНК).
- $\sum_{i=1}^n \hat{u}_i x_i = 0$ (условие ортогональности в МНК).

Теперь раскроем перекрёстное произведение:

$$\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = \sum_{i=1}^n \hat{u}_i (\hat{\beta}_1 + \hat{\beta}_2 x_i - \bar{y}) =$$

(подставим $\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}$):

$$\begin{aligned} &= \sum_{i=1}^n \hat{u}_i (\bar{y} - \hat{\beta}_2 \bar{x} + \hat{\beta}_2 x_i - \bar{y}) = \sum_{i=1}^n \hat{u}_i \hat{\beta}_2 (x_i - \bar{x}) = \hat{\beta}_2 \sum_{i=1}^n \hat{u}_i (x_i - \bar{x}) = \\ &= \hat{\beta}_2 \left(\sum_{i=1}^n \hat{u}_i x_i - \bar{x} \sum_{i=1}^n \hat{u}_i \right) = \hat{\beta}_2 (0 - \bar{x} \cdot 0) = 0. \end{aligned}$$

Таким образом, перекрёстное произведение равно нулю, и равенство $TSS = ESS + RSS$ доказано.

Задача 11. Для модели парной линейной регрессии без константы $y_i = \beta_1 x_i + u_i$ покажите, что в общем случае

- $\sum_{i=1}^n \hat{u}_i \neq 0$
- $\bar{y} \neq \hat{\beta}_1 \bar{x}$
- $TSS \neq RSS + ESS$
- $R^2 = 1 - \frac{RSS}{TSS} \notin [0, 1]$

Решение. а) В модели с константой выполняется $\sum \hat{u}_i = 0$, но в модели **без константы**:

$$\sum_{i=1}^n \hat{u}_i = \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i.$$

Эта разность равна нулю только если $\sum y_i = \hat{\beta}_1 \sum x_i$, что выполняется не всегда.

Контрпример: Пусть $x = (1, 2)^T$, $y = (3, 5)^T$. Тогда:

$$\hat{\beta}_1 = \frac{1 \cdot 3 + 2 \cdot 5}{1^2 + 2^2} = \frac{13}{5} = 2.6$$

Остатки:

$$\hat{u}_1 = 3 - 2.6 \cdot 1 = 0.4$$

$$\hat{u}_2 = 5 - 2.6 \cdot 2 = -0.2$$

Сумма остатков: $0.4 - 0.2 = 0.2 \neq 0$.

б) В модели **без константы**:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i,$$

$$\hat{\beta}_1 \bar{x} = \hat{\beta}_1 \cdot \frac{1}{n} \sum_{i=1}^n x_i.$$

Из пункта «а» следует, что в общем случае $\sum y_i \neq \hat{\beta}_1 \sum x_i$, значит:

$$\bar{y} \neq \hat{\beta}_1 \bar{x}.$$

в) В модели **с константой** выполняется разложение:

$$\text{TSS} = \text{ESS} + \text{RSS},$$

где:

- $\text{TSS} = \sum (y_i - \bar{y})^2$ (общая сумма квадратов),
- $\text{ESS} = \sum (\hat{y}_i - \bar{y})^2$ (объясненная сумма квадратов),
- $\text{RSS} = \sum \hat{u}_i^2$ (остаточная сумма квадратов).

В модели **без константы** это разложение **не выполняется**, так как:

$$\sum \hat{u}_i \neq 0 \Rightarrow \text{выборочная ковариация между } \hat{y}_i \text{ и } \hat{u}_i \text{ не равна нулю.}$$

Следовательно:

$$\text{TSS} \neq \text{ESS} + \text{RSS}.$$

г) Коэффициент детерминации R^2 вычисляется по формуле:

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}.$$

Поскольку для модели без константы $\text{TSS} \neq \text{ESS} + \text{RSS}$, то R^2 может выходить за пределы $[0, 1]$ если $\text{RSS} > \text{TSS}$, то $R^2 < 0$.

Задача 12. Покажите, что для регрессии вида $y_i = \beta_1 + \beta_2 x_i + u_i$ R^2 обладает следующими свойствами:

- $R^2 = \text{sCorr}^2(y, x)$, где $\text{sCorr}^2(y, x)$ – выборочный коэффициент корреляции.
- R^2 для заданной регрессии и R^2 для регрессии $x_i = \alpha_1 + \alpha_2 y_i + v_i$ равны.
- Если $\hat{\beta}_1 = 0$, то $R^2 = 0$.

Задача 13. Пусть мы оценили модель $y_i = \beta_1 + \beta_2 x_i + u_i$ с помощью МНК.

- Если данные оказались центрированными, что вы можете сказать о $\hat{\beta}_1$?
- Ко всем наблюдениям x_i прибавили 15. Как изменятся $\hat{\beta}_1, \hat{\beta}_2$?
- Все наблюдения x_i увеличили в 5 раз, что произойдёт с $\hat{\beta}_1, \hat{\beta}_2$ и \hat{y}_i ?

Решение. а) Если данные центрированы ($\bar{x} = 0, \bar{y} = 0$), то:

$$\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x} = 0 - \hat{\beta}_2 \cdot 0 = 0$$

б) Ко всем x_i прибавили 15. Обозначим новый регрессор как $x'_i = x_i + 15$. Тогда новое среднее равно $\bar{x}' = \bar{x} + 15$.

- Оценка наклона не изменится:

$$\hat{\beta}'_2 = \frac{\sum (x'_i - \bar{x}')(y_i - \bar{y})}{\sum (x'_i - \bar{x}')^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \hat{\beta}_2.$$

- Оценка константы изменится:

$$\hat{\beta}'_1 = \bar{y} - \hat{\beta}'_2 \bar{x}' = \bar{y} - \hat{\beta}_2 (\bar{x} + 15) = \hat{\beta}_1 - 15\hat{\beta}_2.$$

в) Все x_i увеличили в 5 раз. Обозначим новый регрессор как $x''_i = 5x_i$. Тогда новое среднее равно $\bar{x}'' = 5\bar{x}$.

- Оценка наклона изменится:

$$\hat{\beta}''_2 = \frac{\sum (x''_i - \bar{x}'')(y_i - \bar{y})}{\sum (x''_i - \bar{x}'')^2} = \frac{5 \sum (x_i - \bar{x})(y_i - \bar{y})}{25 \sum (x_i - \bar{x})^2} = \frac{\hat{\beta}_2}{5}.$$

- Оценка константы не изменится:

$$\hat{\beta}''_1 = \bar{y} - \hat{\beta}''_2 \bar{x}'' = \bar{y} - \frac{\hat{\beta}_2}{5} \cdot 5\bar{x} = \bar{y} - \hat{\beta}_2 \bar{x} = \hat{\beta}_1.$$

- Прогнозные значения:

$$\hat{y}''_i = \hat{\beta}''_1 + \hat{\beta}''_2 x''_i = \hat{\beta}_1 + \frac{\hat{\beta}_2}{5} \cdot 5x_i = \hat{\beta}_1 + \hat{\beta}_2 x_i = \hat{y}_i.$$

Задача 14. Рассмотрим регрессию $\hat{y}_i = \hat{\beta}_1 z_i + \hat{\beta}_2 x_i$. Все исходные данные поместим в матрицу X и вектор y :

$$X = \begin{pmatrix} z_1 & x_1 \\ \vdots & \vdots \\ z_n & x_n \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

- а) Выпишите явно матрицы X^T , $X^T y$, $X^T X$, $y^T X$, $y^T y$ и укажите их размер.
- б) Выпишите условия первого порядка для оценок $\hat{\beta}_1$ и $\hat{\beta}_2$ по методу наименьших квадратов.
- в) Запишите эти же условия в виде линейной системы

$$\begin{cases} \hat{\beta}_1 \cdot \dots + \hat{\beta}_2 \cdot \dots = \dots \\ \hat{\beta}_1 \cdot \dots + \hat{\beta}_2 \cdot \dots = \dots \end{cases}$$

- г) Как упростится данная система для регрессии $\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_i$?
-

д) Запишите систему условий первого порядка с помощью матрицы X и вектора y ;

Задача 15. Рассмотрим модель $y_i = \beta_1 + \beta_2 x_i + u_i$, где

$$x = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}, \quad y = \begin{pmatrix} 1 \\ 1 \\ 4 \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}, \quad u = \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix}$$

- а) Укажите число наблюдений
- б) Найдите $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$.
- в) С помощью МНК найдите оценку для вектора неизвестных коэффициентов без использования матриц.
- г) Перепишите модель в матричном виде и получите оценку коэффициентов через матричные формулы для МНК (убедитесь, что оценки совпали :)).
- д) Найдите $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$.
- е) Чему равен R^2 в модели? Прокомментируйте полученное значение с точки зрения качества оценённого уравнения регрессии.

Решение. а) Количество строк в векторе y равно 3. Следовательно, $n = 3$.

б) Вычислим среднее значение \bar{y} :

$$\bar{y} = \frac{1 + 1 + 4}{3} = 2.$$

Тогда общая сумма квадратов:

$$TSS = (1 - 2)^2 + (1 - 2)^2 + (4 - 2)^2 = 1 + 1 + 4 = 6.$$

в) Формулы для МНК-оценок:

$$\hat{\beta}_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2},$$
$$\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}.$$

Вычислим необходимые величины:

$$\bar{x} = \frac{1 + 2 + 3}{3} = 2,$$

$$\sum (x_i - \bar{x})(y_i - \bar{y}) = (1 - 2)(1 - 2) + (2 - 2)(1 - 2) + (3 - 2)(4 - 2) = 1 + 0 + 2 = 3,$$

$$\sum (x_i - \bar{x})^2 = (1 - 2)^2 + (2 - 2)^2 + (3 - 2)^2 = 1 + 0 + 1 = 2.$$

Тогда оценки:

$$\hat{\beta}_2 = \frac{3}{2} = 1.5,$$

$$\hat{\beta}_1 = 2 - 1.5 \times 2 = -1.$$

г) Матрица регрессоров и вектор зависимой переменной:

$$X = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{pmatrix}, \quad Y = \begin{pmatrix} 1 \\ 1 \\ 4 \end{pmatrix}.$$

МНК-оценка коэффициентов:

$$\hat{\beta} = (X^T X)^{-1} X^T Y.$$

Вычислим:

$$X^T X = \begin{pmatrix} 3 & 6 \\ 6 & 14 \end{pmatrix}, \quad X^T Y = \begin{pmatrix} 6 \\ 15 \end{pmatrix},$$

$$(X^T X)^{-1} = \frac{1}{6} \begin{pmatrix} 14 & -6 \\ -6 & 3 \end{pmatrix},$$

$$\hat{\beta} = \frac{1}{6} \begin{pmatrix} 14 \times 6 - 6 \times 15 \\ -6 \times 6 + 3 \times 15 \end{pmatrix} = \begin{pmatrix} -1 \\ 1.5 \end{pmatrix}.$$

Результаты совпали с предыдущим пунктом.

д) Прогнозные значения:

$$\hat{y}_1 = -1 + 1.5 \times 1 = 0.5,$$

$$\hat{y}_2 = -1 + 1.5 \times 2 = 2,$$

$$\hat{y}_3 = -1 + 1.5 \times 3 = 3.5.$$

Остаточная сумма квадратов:

$$\text{RSS} = (1 - 0.5)^2 + (1 - 2)^2 + (4 - 3.5)^2 = 0.25 + 1 + 0.25 = 1.5.$$

е) Формула для R^2 :

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{1.5}{6} = 0.75.$$

Интерпретация: модель объясняет 75% вариации зависимой переменной, что говорит о хорошем качестве подгонки модели.

Задача 16. Рассмотрим модель $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + u_i$, где

$$X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}, \quad y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}, \quad u = \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \end{pmatrix}$$

Для удобства расчётов даны матрицы:

$$X^T X = \begin{pmatrix} 5 & 2 & 1 \\ 2 & 2 & 1 \\ 1 & 1 & 1 \end{pmatrix}, \quad (X^T X)^{-1} = \begin{pmatrix} 1/3 & -1/3 & 0 \\ -1/3 & 4/3 & -1 \\ 0 & -1 & 2 \end{pmatrix}.$$

- а) Укажите число наблюдений
- б) Укажите число регрессоров в модели, учитывая свободный член.
- в) Найдите $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$.
- г) С помощью МНК найдите оценку для вектора неизвестных коэффициентов.
- д) Найдите вектор прогнозов \hat{y} .
- е) Найдите $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$.
- ж) Чему равен R^2 в модели? Прокомментируйте полученное значение с точки зрения качества оценённого уравнения регрессии.

Решение. а) Количество строк в матрице X равно 5, следовательно, $n = 5$.

б) В модель включена свободный член и две объясняющие переменные (x и z), то есть $k = 3$.

в) Вычислим среднее значение \bar{y} :

$$\bar{y} = \frac{1 + 2 + 3 + 4 + 5}{5} = 3.$$

Тогда общая сумма квадратов:

$$TSS = (1 - 3)^2 + (2 - 3)^2 + (3 - 3)^2 + (4 - 3)^2 + (5 - 3)^2 = 4 + 1 + 0 + 1 + 4 = 10.$$

г) Используем матричную формулу МНК:

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

Умножим обратную матрицу на $X^T y$:

$$X^T y = \begin{pmatrix} 1 + 2 + 3 + 4 + 5 \\ 0 + 0 + 0 + 4 + 5 \\ 0 + 0 + 0 + 0 + 5 \end{pmatrix} = \begin{pmatrix} 15 \\ 9 \\ 5 \end{pmatrix},$$

$$(X^T X)^{-1} X^T y = \begin{pmatrix} 1/3 & -1/3 & 0 \\ -1/3 & 4/3 & -1 \\ 0 & -1 & 2 \end{pmatrix} \begin{pmatrix} 15 \\ 9 \\ 5 \end{pmatrix} = \begin{pmatrix} 5 - 3 + 0 \\ -5 + 12 - 5 \\ 0 - 9 + 10 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix}.$$

д) Вычислим вектор прогнозов $\hat{y} = X\hat{\beta}$:

$$\hat{y} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \\ 2 \\ 4 \\ 5 \end{pmatrix}.$$

е) Остаточная сумма квадратов:

$$RSS = (1 - 2)^2 + (2 - 2)^2 + (3 - 2)^2 + (4 - 4)^2 + (5 - 5)^2 = 1 + 0 + 1 + 0 + 0 = 2.$$

ж) Формула для R^2 :

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{2}{10} = 0.8.$$

Задача 17. Константин оценивает влияние продаж на заработную плату менеджера: $salary_i = \beta_1 + \beta_2 sales_i + u_i$. Он оценил коэффициенты с помощью МНК, нашёл R_1^2 и ему не понравился результат.

- а) Тогда Константин выкинул одно из наблюдений, переоценил модель и получил R_2^2 . Может ли он сравнить модели по R^2 и выбрать наилучшую?
- б) Константин добавил в модель регрессор числа созвонов с начальством $sinks_i$. Покажите, что R^2 у данной регрессии вырастет в любом случае.
- в) Константин слышал, что некоторые исследователи логарифмируют зарплаты перед включением в модель. Объясните, можно ли сравнивать по R^2 регрессии $salary_i$ и $\log(salary_i)$ на один и тот же набор регрессоров?

Задача 18. Я очень хочу тут реальный датасет с точными датами рождения людей :)

Задача 19. Рассмотрим модель множественной регрессии

$$y = X\beta + u,$$

где X — матрица признаков размерности $[n \times k]$. Определим матрицу $M = I - X(X^T X)^{-1} X^T$. Покажите, что

- а) матрица M симметричная;
- б) матрица M идемпотентная;
- в) $\text{trace } M = n - k$.

Решение. а) Вычислим транспонированную матрицу M :

$$M^T = (I_n - X(X^T X)^{-1} X^T)^T = I_n^T - (X(X^T X)^{-1} X^T)^T = I_n - X((X^T X)^{-1})^T X^T.$$

Поскольку $X^T X$ симметрична, то $(X^T X)^{-1}$ также симметрична:

$$((X^T X)^{-1})^T = (X^T X)^{-1}.$$

Следовательно,

$$M^T = I_n - X(X^T X)^{-1} X^T = M.$$

Таким образом, M симметрична.

- б) Матрица M называется идемпотентной, если $M^2 = M$.

Вычислим M^2 :

$$M^2 = (I_n - X(X^T X)^{-1} X^T) (I_n - X(X^T X)^{-1} X^T) \quad (1)$$

$$= I_n - 2X(X^T X)^{-1} X^T + X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T. \quad (2)$$

Упростим последнее слагаемое:

$$X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T = X(X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T.$$

Подставим обратно:

$$M^2 = I_n - 2X(X^T X)^{-1} X^T + X(X^T X)^{-1} X^T = I_n - X(X^T X)^{-1} X^T = M.$$

Таким образом, M идемпотентна.

в) Используем свойства следа матрицы:

$$\text{trace}(M) = \text{trace}(I_n - X(X^T X)^{-1} X^T) \quad (3)$$

$$= \text{trace}(I_n) - \text{trace}(X(X^T X)^{-1} X^T). \quad (4)$$

След единичной матрицы равен её размеру, $\text{trace}(I_n) = n$. Для второго слагаемого воспользуемся свойством следа $\text{trace}(AB) = \text{trace}(BA)$:

$$\text{trace}(X(X^T X)^{-1} X^T) = \text{trace}((X^T X)^{-1} X^T X) = \text{trace}(I_k) = k.$$

Таким образом,

$$\text{trace}(M) = n - k.$$

Задача 20. Пусть $s = (1 \ 1 \ \dots \ 1)^T$ — вектор размерности $[n \times 1]$, состоящий из единиц. Определим матрицу $\pi = s^T (s^T s)^{-1} s^T$. Матрица π — это матрица размерности $[n \times n]$ вида

$$\pi = \frac{1}{n} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \dots & \dots & \dots & \dots \\ 1 & 1 & \dots & 1 \end{pmatrix}.$$

Покажите, что $\pi \cdot c = \bar{c} \cdot s$, где c — произвольный вектор размерности $[n \times 1]$; \bar{c} — среднее арифметическое, посчитанное по элементам вектора c .

Решение. Вычислим произведение $\pi \cdot c$:

$$\pi \cdot c = \frac{1}{n} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{pmatrix}$$

Умножение матрицы на вектор дает:

$$\pi \cdot c = \frac{1}{n} \begin{pmatrix} \sum_{i=1}^n c_i \\ \sum_{i=1}^n c_i \\ \vdots \\ \sum_{i=1}^n c_i \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n c_i \cdot \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} \bar{c} \\ \bar{c} \\ \vdots \\ \bar{c} \end{pmatrix}.$$

По определению среднего значения:

$$\pi \cdot c = \bar{c} \cdot s,$$

где $s = (1, 1, \dots, 1)^T$.

Таким образом, мы показали, что умножение матрицы π на произвольный вектор c дает вектор, все элементы которого равны среднему значению элементов вектора c .

6. Предпосылки о математическом ожидании и дисперсии

В этой главе мы познакомимся с понятиями независимости и линейной независимости, расчётом математических ожиданий, ковариаций и дисперсий в матричном виде.

Добавим в метод наименьших квадратов ряд статистических предпосылок на ожидание и дисперсию.

Сформулируем и докажем теорему Гаусса - Маркова (которая пообещает, что МНК-оценки будут обладать свойствами несмещённости и эффективности).

6.1. Иерархия зависимостей случайных величин

Напомним определение наилучшей линейной аппроксимации.

Определение 6.1 (наилучшая линейная аппроксимация). Наилучшее линейное приближение величины r с помощью величины s — это линейная функция от s ,

$$\text{BestLin}(r | s) = \beta_1 + \beta_2 s,$$

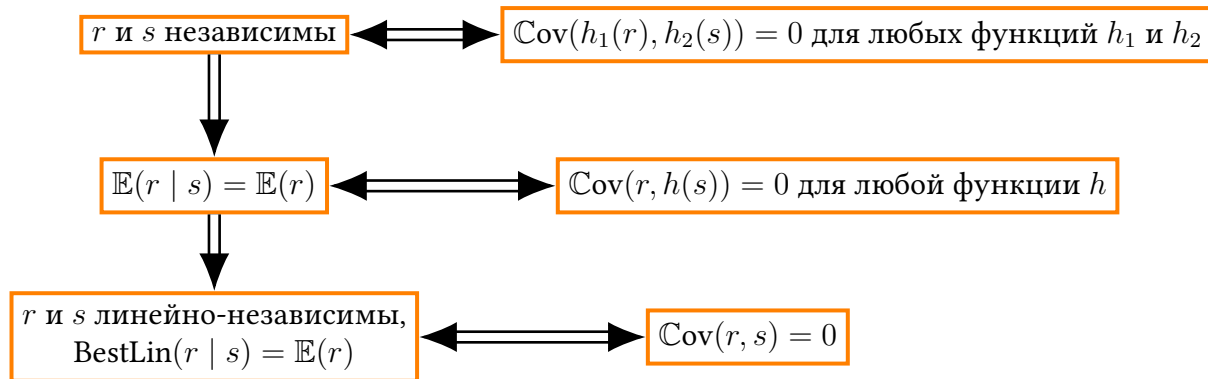
где константы β_1 и β_2 находятся из решения задачи оптимизации $\mathbb{E}((r - \text{BestLin}(r, s))^2) \rightarrow \min_{\beta_1, \beta_2}$. При решении задачи оказывается, что

$$\beta_1 = \mathbb{E}(r) - \frac{\text{Cov}(r, s)}{\text{Var}(s)} \mathbb{E}(s), \quad \beta_2 = \frac{\text{Cov}(r, s)}{\text{Var}(s)}.$$

Определение 6.2 (линейная независимость). Величины r и s называются линейно-независимыми, если $\text{BestLin}(r | s) = \mathbb{E}(r)$.

Некоторые авторы считают условие $\text{Cov}(r, s) = 0$ определением линейной независимости.

Можно выделить три степени независимости случайных величин. Рассмотрим их на примере пары произвольных величин r и s .



Далее напомним определение независимых случайных величин.

Определение 6.3 (независимость случайных величин). Случайные величины r и s называются независимыми если для любых⁴ числовых множеств A и B независимы события $\{r \in A\}$ и $\{s \in B\}$:

$$\mathbb{P}(r \in A, s \in B) = \mathbb{P}(r \in A) \cdot \mathbb{P}(s \in B)$$

⁴Не совсем любых, требуется измеримость множеств. В рамках нашего курса мы не будем обращать внимания на данный нюанс.

Из независимости величин r и s следует, что информация, известная об s , никак не помогает угадывать значение r . Поэтому условное математическое ожидание для r равно безусловному. Точно также из независимости r и s следует $\mathbb{E}(s \mid r) = \mathbb{E}(s)$. Обратное утверждение неверно, что показывает контрпример ниже.

Задача 21. Покажем, что из равенства условного и безусловного математических ожиданий не следует независимость случайных величин. Пусть дискретные случайные величины r характеризуют погоду (-1 снег, 1 солнце, 0 дождь), s — наличие зонта (0 нет или 1 есть) и ниже приведена таблица их совместного распределения.

	1/3	1/3	1/3
r	-1	1	0
s	0	0	1

Решение. Уже по формулировке подозреваем, что величины зависимые :).

Найдём условное ожидание зонта при условии, что мы видим погоду на улице: $\mathbb{E}(s \mid r) = \begin{cases} 0 & \text{если } r \in \{-1, 1\} \\ 1, & \text{если } r = 0. \end{cases}$

Получается, что информация о погоде помогает предсказать наличие зонтика, события не являются независимыми.

Найдём ожидания о погоде за окном, если вы можете наблюдать наличие или отсутствие зонта у человека: $\mathbb{E}(r \mid s) = \begin{cases} (-1) \times 1/6 + 1 \times 1/6 = 0, & \text{если } s = 0, \\ 0, & \text{если } s = 1. \end{cases}$

Обычное безусловное ожидание погоды на улице: $\mathbb{E}(r) = (-1) \times 1/3 + 1 \times 1/3 + 0 \times 1/3 = 0$

Получается, что $\mathbb{E}(r \mid s) = \mathbb{E}(r) = 0$, но события зависимы.

Вернемся к тому факту, что из равенства условного и безусловного математических ожиданий следует нулевая ковариация. Используя закон повторных математических ожиданий $\text{Cov}(r, s) = \mathbb{E}(rs) - \mathbb{E}(r) \mathbb{E}(s) = \mathbb{E}(\mathbb{E}(rs \mid s)) - \mathbb{E}(r) \mathbb{E}(s) = \mathbb{E}(s \mathbb{E}(r \mid s)) - \mathbb{E}(r) \mathbb{E}(s) = \mathbb{E}(r) \mathbb{E}(s) - \mathbb{E}(r) \mathbb{E}(s) = 0$.

Задача 22. Из нулевой ковариации не следует равенство условного и безусловного математических ожиданий (и тем более не следует независимость). Пусть случайная величина s имеет равномерное распределение на отрезке $[-1; 1]$, а $r = s^2$.

Решение. Напоминаем, что для равномерно распределённой случайной величины $\mathbb{E}(s) = \frac{-1+1}{2} = 0$, $pdf(s) = \frac{1}{1-(-1)} = \frac{1}{2}$.

$\mathbb{E}(r \mid s) = \mathbb{E}(s^2 \mid s) = s^2 \neq 0$ в общем случае.

При этом $\text{Cov}(r, s) = \mathbb{E}(rs) - \mathbb{E}(r) \mathbb{E}(s) = \mathbb{E}(s^3) - \mathbb{E}(s^2) \times 0 = \mathbb{E}(s^3)$.

Математическое ожидание сложной функции $\mathbb{E}(g(x)) = \int_b^r g(x) pdf(x) dx$, если $x \in [a, b]$.

Найдём $\mathbb{E}(s^3) = \int_{-1}^1 s^3 pdf(s) ds = \int_{-1}^1 s^3 \frac{1}{2} ds = \frac{1}{8} s^4 \Big|_{-1}^1 = 0$. Значит, мы получили нулевую ковариацию у зависимых случайных величин.

Вывод

Существуют независимые случайные величины, но на ???

6.2. Ожидание и ковариационная матрица

Любопытный читатель снова может заглянуть в «Поваренную книга любителя матриц» Петерсона и Педерсона. [PP12].

Пусть r — случайный вектор размерности $[n \times 1]$, s — случайный вектор размерности $[k \times 1]$, A и b — неслучайные матрица и вектор соответственно, имеющие подходящие размерности.

Математическим ожиданием случайного вектора r называется вектор

$$\mathbb{E}(r) = \begin{pmatrix} \mathbb{E}(r_1) \\ \mathbb{E}(r_2) \\ \dots \\ \mathbb{E}(r_n) \end{pmatrix}.$$

Ковариационная матрица вектора r определяется следующим образом:

$$\mathbb{V}\text{ar}(r) = \begin{pmatrix} \text{Cov}(r_1, r_1) & \text{Cov}(r_1, r_2) & \dots & \text{Cov}(r_1, r_n) \\ \text{Cov}(r_2, r_1) & \text{Cov}(r_2, r_2) & \dots & \text{Cov}(r_2, r_n) \\ \dots & \dots & \dots & \dots \\ \text{Cov}(r_n, r_1) & \text{Cov}(r_n, r_2) & \dots & \text{Cov}(r_n, r_n) \end{pmatrix}.$$

Ковариационная матрица векторов r и s определяется следующим образом:

$$\text{Cov}(r, s) = \begin{pmatrix} \text{Cov}(r_1, s_1) & \text{Cov}(r_1, s_2) & \dots & \text{Cov}(r_1, s_k) \\ \text{Cov}(r_2, s_1) & \text{Cov}(r_2, s_2) & \dots & \text{Cov}(r_2, s_k) \\ \dots & \dots & \dots & \dots \\ \text{Cov}(r_n, s_1) & \text{Cov}(r_n, s_2) & \dots & \text{Cov}(r_n, s_k) \end{pmatrix}.$$

Свойства вектора математических ожиданий и ковариационной матрицы:

- а) $\mathbb{E}(Ar + b) = A \mathbb{E}(r) + b$
- б) $\text{Cov}(r, s) = \mathbb{E}(rs^T) - \mathbb{E}(r) \mathbb{E}(s^T)$
- в) $\text{Cov}(Ar + b, s) = A \text{Cov}(r, s)$
- г) $\text{Cov}(r, As + b) = \text{Cov}(r, s) A^T$
- д) $\mathbb{V}\text{ar}(r) = \text{Cov}(r, r) = \mathbb{E}(rr^T) - \mathbb{E}(r) \mathbb{E}(r^T)$
- е) $\mathbb{V}\text{ar}(Ar + b) = A \mathbb{V}\text{ar}(r) A^T$
- ж) $\mathbb{E}(r^T Ar) = \text{trace}(A \mathbb{V}\text{ar}(r)) + \mathbb{E}(r^T) A \mathbb{E}(r)$
- з) Если вектора r и s имеют одинаковый размер, то $\mathbb{V}\text{ar}(r + s) = \mathbb{V}\text{ar}(r) + \mathbb{V}\text{ar}(s) + \text{Cov}(r, s) + \text{Cov}(s, r)$

Условные ожидание и дисперсия определяются аналогично и обладают аналогичными свойствами. Главное — не забывать ставить вертикальную палочку!

Например,

написать пример

6.3. Теорема Гаусса — Маркова для парной линейной регрессии

Вернёмся к модели парной линейной регрессии, для которой мы научились получать оценки $\hat{\beta}_1, \hat{\beta}_2$ методом наименьших квадратов:

$$y_i = \beta_1 + \beta_2 x_i + u_i.$$

Мы убедились, что практически для любых данных в y и x оценки технически могут быть получены (кроме случая, когда все x_i являются одинаковыми и в знаменателе для $\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$ произойдёт деление на ноль).

При этом мы ничего не говорили о качестве полученных оценок, а оказывается, что МНК-оценки далеко не всегда будут "хорошими" с точки зрения несмещённости, состоятельности и эффективности.

Сформулируем теорему, которая будет отвечать на вопрос, в каких случаях можно доверять оценкам, полученным методом наименьших квадратов.

Теорема 6.4 (Гаусс — Марков для парной линейной регрессии). Если

1. **Линейность:** модель линейна по параметрам: $y_i = \beta_1 + \beta_2 x_i + u_i$;
2. **Экзогенность:** условное ожидание случайных ошибок равно нулю: $\mathbb{E}(u_i \mid x_1, x_2, \dots, x_n) = 0$;
3. **Гомоскедастичность:** условная дисперсия случайных ошибок постоянна: $\text{Var}(u_i \mid x_1, \dots, x_n) = \sigma_u^2, \quad 0 < \sigma_u^2 < \infty$;
4. **Линейная независимость:** $\mathbb{E}(u_i u_j \mid x_1, x_2, \dots, x_n) = 0$;
5. Оценки $\hat{\beta}_1, \hat{\beta}_2$ получены методом наименьших квадратов;

то

- (a) Оценки $\hat{\beta}_1, \hat{\beta}_2$ являются линейными по y ;
- (b) Оценки $\hat{\beta}_1, \hat{\beta}_2$ являются условно несмещёнными, $\mathbb{E}(\hat{\beta}_i \mid x_1, \dots, x_n) = \beta_i$ и несмещёнными, $\mathbb{E}(\hat{\beta}_i) = \beta_i$;
- (c) Оценка любого коэффициента $\hat{\beta}_i$ является наиболее эффективной в классе линейных несмещённых оценок.

куда-то присобачить состоятельность, тогда появятся конечные четвёртые моменты? $E[X^4] < \infty, \quad E[Y^4] < \infty$

Перед тем, как перейти к доказательству, обсудим, что означает каждая предпосылка и какие ограничения она накладывает.

Линейность модели

В данном случае подразумевается, что модель является линейной по параметрам β_1, β_2 . Это является достаточно сильным упрощением реальности. Мы не можем быть уверенными, что например, зарплата y линейно возрастает по опыту x .

Если на самом деле окажется, что мы пытаемся оценить нелинейные эффекты с помощью линейной модели, то качество оценок будет печальным и никакая теорема Гаусса-Маркова уже не поможет. Мы в дальнейшем ещё будем рассуждать о том, как заподозрить нелинейность в ваших данных.

В качестве утешения уже сейчас можно строить нелинейные по x модели: например $y_i = \beta_1 + \beta_2 x_i^2 + u_i$ удовлетворяет условию теоремы Гаусса - Маркова. Модель $y_i = \beta_1 + e^{\beta_2 x_i} + u_i$ не удовлетворяет условию, так как нелинейна по β_2 .

Экзогенность

Нулевое условное математическое ожидание (или экзогенность) выписывается в форме

$$E(u_i | x_1, x_2, \dots, x_n) = 0.$$

связка с другим разделом

Мы предполагаем, что в имеющихся данных по x_i не осталось никакой информации, которая могла бы выдать полезный, отличный от нуля прогноз для u_i . По предыдущим разделам мы можем выписать более слабую версию этой предпосылки: $E[u_i | x_i] = 0$ и даже ещё более слабые $E(u_i) = 0$ и $\text{Cov}(x, u) = 0$. В такой форме чуть легче понять, что в реальной жизни можно считать нарушением предпосылки об экзогенности.

Так как $E(u_i | x_i) = 0 \Rightarrow \text{Cov}(x, u) = 0$, значит по правилам построения отрицаний $\text{Cov}(x, u) \neq 0 \Rightarrow E(u_i | x_i) \neq 0$. Если получится придумать, по какой причине u в вашей модели коррелирует с x , вы точно нарушаете предпосылку. Например, вы всё ещё оцениваете влияние опыта работы x на зарплату y . Все факторы, влияющие на заработную плату, которые не были учтены в модели, автоматически перемещаются в ошибки (возраст, пол, интеллект, стрессоустойчивость, коммуникабельность и т.д.). Тот же возраст часто коррелирует с опытом работы x , поэтому в такой парной линейной регрессии нарушается предпосылка об экзогенности (или, иными словами, присутствует эндогенность).

Нарушение предпосылки об экзогенности одно из самых распространённых. Всегда можно придумать регрессоры, которые забыли или не смогли включить в модель, а они коррелируют с уже имеющимися регрессорами. Мы отдельно вернёмся к этой теме в следующих разделах.

Гомоскедастичность

Условная дисперсия случайной ошибки постоянна (гомоскедастичность):

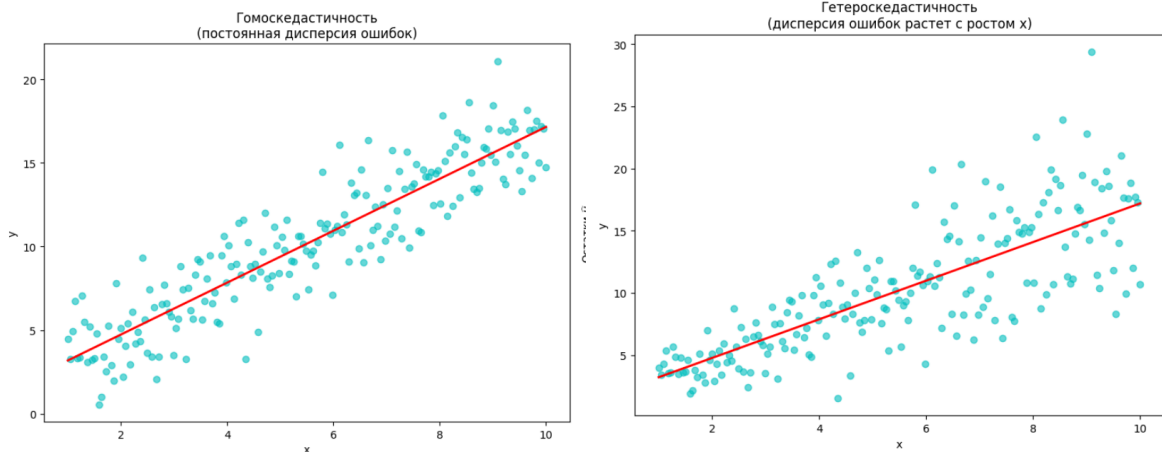
$$\text{Var}(u_i | x_1, \dots, x_n) = \sigma_u^2, \quad 0 < \sigma_u^2 < \infty$$

На рисунке ?? приведены примеры для иллюстрации гомоскедастичности и гетероскедастичности. В случае гетероскедастичности разброс y выше при больших x , что намекает о зависимости дисперсии ошибок от регрессора. В гомоскедастичность иногда непросто поверить: для той же регрессии зарплаты на опыт работы наверняка можно предположить, что при отсутствии опыта сотрудники будут получать похожие офферы. По мере получения опыта одинаково долго проработавшие в компании сотрудники могут уже получать существенно разные зарплаты из-за иных, неучтённых пока в модели факторов. Это приводит к тому, что некоторые исследователи заранее не верят в гомоскедастичность и предпринимают меры по коррекции результатов оценивания.

Во-первых, где в нашем МНК-оценивании используется гомоскедастичность в парной линейной регрессии $y_i = \beta_1 + \beta_2 x_i + u_i$?

Доказательство. Мы уже знаем оценку углового коэффициента

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum y_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = \frac{\sum (\beta_0 + \beta_1 x_i + u_i)(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = \\ &= \frac{\beta_0 \sum (x_i - \bar{x}) + \beta_1 \sum x_i (x_i - \bar{x}) + \sum u_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = \beta_1 + \frac{\sum u_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \end{aligned}$$



Дисперсия оценки $\hat{\beta}_1$:

$$\text{Var}(\hat{\beta}_1 \mid x_1, x_2, \dots, x_n) = \text{Var} \left(\frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2} \mid x_1, x_2, \dots, x_n \right) = \frac{1}{(\sum (x_i - \bar{x})^2)^2} \text{Var} \left(\sum (x_i - \bar{x}) u_i \mid x_1, x_2, \dots, x_n \right)$$

Так как наблюдения независимы и случайные ошибки независимы, дисперсия суммы равна сумме дисперсий:

$$\text{Var}(\hat{\beta}_1 \mid x_1, x_2, \dots, x_n) = \frac{\sum (x_i - \bar{x})^2 \text{Var}(u_i \mid x_1, x_2, \dots, x_n)}{(\sum (x_i - \bar{x})^2)^2} = \frac{\sigma_u^2}{\sum (x_i - \bar{x})^2}.$$

□

Такую красивую короткую формулу для теоретической дисперсии оценки углового коэффициента мы получаем только в предпосылке о постоянстве дисперсии случайной ошибки. Без неё пришлось бы делать предположение о том, как ведёт себя ошибка для каждого значений регрессора (чем мы тоже будем заниматься далее).

Что будет, если предположить гомоскедастичность там, где её нет? Показано, что расчёты **в условиях гомоскедастичности обычно недооценивают** дисперсию оценок коэффициентов. Это приводит к завышенным стандартным ошибкам и проблемам в тестировании гипотез:

надо ссылку пруф

Например, для тестирования $H_0 : \beta_1 = 0$ против $H_1 : \beta_1 \neq 0$ с использованием t-статистики мы получим

$$t_{\text{stat}} = \frac{\hat{\beta}_1 - 0}{s.e.(\hat{\beta}_1)}$$

Заниженная стандартная ошибка $s.e.(\hat{\beta}_1)$ приведёт к слишком большому значению статистики и, значит, к большей вероятности отвергнуть гипотезу о незначимости коэффициента. Получается, что расчёты при гомоскедастичности могут ошибочно говорить о значимости регрессора. В разделе о гетероскедастичности мы будем говорить о том, как можно исправить ситуацию и проводить более пессимистичное тестирование гипотез.

Линейная независимость

Под линейной независимостью подразумевается, что

$$\mathbb{E}(u_i u_j \mid x_1, x_2, \dots, x_n) = 0$$

В более слабой форме в предыдущем разделе это можно записать как $\text{Cov}(u_i, u_j \mid x_1, x_2, \dots, x_n) = 0$ $\forall i \neq j$. Это свойство можно считать автоматически выполненным, если перед нами независимая выборка наблюдений. Тогда все неучтённые факторы, влияющие на зарплату человека i , не связаны с неучтёнными факторами человека j .

Когда эта предпосылка не выполняется: например, выборка состоит из сотрудников одной компании, в которой опытные сотрудники сговорились и забирают себе все премии, а неопытные ничего не получают. Получается, что в необъяснённых факторах для зарплаты будет корреляция определённого знака в зависимости от опыта работы.

ну не самый удачный пример?

Теперь, когда мы поняли, какой смысл кроется за предпосылками об экзогенности, гомоскедастичности и линейной независимости наблюдений, перейдем к следствиям теоремы. Оценки должны получаться несмещённые и наиболее эффективные в классе линейных.

В иностранной литературе для простоты запоминания используется аббревиатура BLUE, best linear unbiased estimator. Хорошие оценки подобны хорошему подвенечному платью,

Something Olde, Something New, Something Borrowed, Something Blue, A Sixpence in your Shoe.

Несмещённость

Если предпосылка об экзогенности $\mathbb{E}(u_i \mid x_1, x_2, \dots, x_n) = 0$ выполняется, то оценка получается **несмещённая**. Можно надеяться на то, что она попала недалеко от истинного значения параметра и поэтому имеет какой-то физический смысл (а значит, её можно **интерпретировать**):

Доказательство. В парной линейной регрессии $y_i = \beta_1 + \beta_2 x_i + u_i$ мы уже знаем оценку углового коэффициента

$$\hat{\beta}_1 = \beta_1 + \frac{\sum u_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

Вспомним определение несмещённости:

$$\mathbb{E}(\hat{\beta}_1 \mid x_1, x_2, \dots, x_n) = \mathbb{E} \left(\beta_1 + \frac{\sum u_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \mid x_1, x_2, \dots, x_n \right) = \beta_1 + \frac{\sum (x_i - \bar{x}) \mathbb{E}(u_i \mid x_1, x_2, \dots, x_n)}{\sum (x_i - \bar{x})^2}$$

При условии $\mathbb{E}(u_i \mid x_1, x_2, \dots, x_n) = 0$, получаем: $\mathbb{E}[\hat{\beta}_1 \mid x_1, x_2, \dots, x_n] = \beta_1$.

По закону повторного математического ожидания $\mathbb{E}(\mathbb{E}[\hat{\beta}_1 \mid x_1, x_2, \dots, x_n]) = \mathbb{E}(\hat{\beta}_1) = \beta_1$ □

Линейность

Оценку $\hat{\beta}_1$ можно представить как линейную комбинацию y_i :

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} = \sum a_i y_i, \text{ где } a_i = \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2}.$$

Таким образом, $\hat{\beta}_1$ – линейная по y оценка.

Эффективность

Под эффективностью мы подразумеваем, что для любой несмещённой и линейной по y оценки $\hat{\beta}_1^{alt}$ выполняется $\text{Var}(\hat{\beta}_1^{alt}) \geq \text{Var}(\hat{\beta}_1)$

Доказательство. Пусть

$$\hat{\beta}_1 = \sum a_i y_i = \sum a_i (\beta_0 + \beta_1 x_i + u_i) = \beta_0 \sum a_i + \beta_1 \sum a_i x_i + \sum a_i u_i.$$

a_i – это какие-то комбинации из x_1, x_2, \dots, x_n , поэтому

$$\mathbb{E}(\hat{\beta}_1 \mid x_1, x_2, \dots, x_n) = \beta_0 \sum a_i + \beta_1 \sum a_i x_i + \sum a_i \mathbb{E}(u_i \mid x_1, x_2, \dots, x_n) = \beta_0 \sum a_i + \beta_1 \sum a_i x_i$$

Для несмещённости $\mathbb{E}(\hat{\beta}_1 \mid x_1, x_2, \dots, x_n) = \beta_1$, значит

$$\begin{cases} \sum a_i = 0, \\ \sum a_i x_i = 1. \end{cases}$$

Альтернативная оценка $\hat{\beta}_1^{alt} = \sum a_i^{alt} y_i$ тоже несмещённая, значит

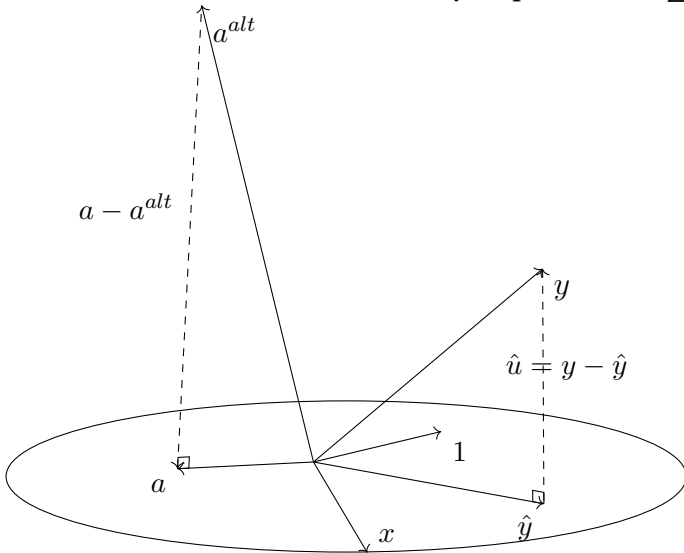
$$\begin{cases} \sum a_i^{alt} = 0, \sum (a_i - a_i^{alt}) = 0 \\ \sum a_i^{alt} x_i = 1, \sum (a_i - a_i^{alt}) x_i = 0. \end{cases}$$

Дисперсия МНК-оценки в условиях гомоскедастичности:

$$\text{Var}(\hat{\beta}_1 \mid x_1, x_2, \dots, x_n) = \text{Var} \left(\sum a_i u_i \mid x_1, x_2, \dots, x_n \right) = \sum a_i^2 \text{Var}(u_i \mid x_1, x_2, \dots, x_n) = \sigma_u^2 \sum a_i^2$$

Аналогично $\text{Var}(\hat{\beta}_1^{alt} \mid x_1, x_2, \dots, x_n) = \sigma_u^2 \sum (a_i^{alt})^2$.

Осталось понять, можем ли мы утверждать, что $\sum (a_i^{alt})^2 \geq \sum a_i^2$?



Оказывается, можем. Вектор a является суммой вектора x и вектора констант (который прячется в \bar{x}). Поэтому a лежит в плоскости регрессоров модели. Судя по $\sum (a_i - a_i^{alt}) = 0, \sum (a_i - a_i^{alt}) x_i = 0$, вектор $a - a^{alt}$ ортогонален вектору констант и вектору x , значит, ортогонален плоскости. Тогда вектор a^{alt} – соединяющая концы двух векторов гипотенуза, длина которой не меньше катетов, значит $\sum (a_i^{alt})^2 \geq \sum a_i^2$.

Если не очень нравится стереометрия, можно воспользоваться раскрытием скобок:

Пусть $\hat{\beta}_1^{alt} = \sum (a_i + d_i)y_i$. Её дисперсия:

$$\mathbb{Var}(\hat{\beta}_1^{alt} \mid x_1, x_2, \dots, x_n) = \sigma_u^2 \sum (a_i + d_i)^2$$

$$\sum (a_i + d_i)^2 = \sum a_i^2 + \sum d_i^2 + 2 \sum a_i d_i = \sum a_i^2 + \sum d_i^2 \geq \sum a_i^2 \quad (\text{т.к. } \sum a_i d_i = 0)$$

$$\sum a_i d_i = \sum \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} d_i = \frac{\sum x_i d_i - \bar{x} \sum d_i}{\sum (x_i - \bar{x})^2} = 0$$

$$\text{т. к. } \sum d_i = \sum (a_i - a_i^{alt}) = 0 \quad \sum x_i d_i = \sum (a_i - a_i^{alt}) x_i = 0$$

□

Состоятельность

Состоятельность оценки: при увеличении объёма выборки оценка сходится по вероятности к истинному значению параметра:

$$\hat{\beta}_n \xrightarrow{P} \beta \quad \text{или} \quad \text{plim}_{n \rightarrow \infty} \hat{\beta}_n = \beta$$

Это означает, что для любого $\varepsilon > 0$:

$$\lim_{n \rightarrow \infty} P(|\hat{\beta}_n - \beta| > \varepsilon) = 0$$

Доказательство. В нашем случае:

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\text{plim } \hat{\beta}_1 = \frac{\mathbb{Cov}(X, Y)}{\mathbb{Var}(X)} = \frac{\mathbb{Cov}(X, \beta_0 + \beta_1 X + u)}{\mathbb{Var}(X)} = \frac{\mathbb{Cov}(X, \beta_1 X) + \mathbb{Cov}(X, u)}{\mathbb{Var}(X)} = \beta_1 + \frac{\mathbb{Cov}(X, u)}{\mathbb{Var}(X)}$$

При выполнении условия экзогенности $\mathbb{Cov}(X, u) = 0$, получаем:

$$\text{plim } \hat{\beta}_1 = \beta_1$$

□

6.4. Теорема Гаусса — Маркова в общем виде

Модель множественной регрессии в матричном виде выглядит следующим образом:

$$y = X\beta + u.$$

Чтобы исследовать свойства полученной точечной оценки $\hat{\beta}$ нам потребуются предпосылки о математическом ожидании и ковариационной матрице вектора u .

Мы предположим, что случайные ошибки в среднем равны нулю, а именно,

$$\mathbb{E}(u \mid X) = 0.$$

Предпосылку о математическом ожидании можно записать и в скалярном виде,

$$\mathbb{E}(u_i | X) = 0, \quad \text{при } \forall i \in \{1, \dots, n\}.$$

Важно пояснить смысл введённой предпосылки. При оценивании связи между регрессорами X и переменной y мы не предполагаем, что величины u и X независимы. В ошибки модели попадают все те факторы, которые мы забыли включить в регрессию. Эти факторы могут быть взаимосвязаны с тем, что в регрессию всё же попало. Мы делаем более слабое предположение лишь о бесполезности всей собранной в X информации для угадывания u (и следующей из неё линейной независимости между ошибками и регрессорами, в том числе о нулевой ковариации).

Теорема 6.5 (Гаусс — Марков). Если

1. Модель линейна по параметрам: $y = X\beta + u$;
2. Матрица X размера $[n \times k]$ имеет полный ранг k .
3. Условное ожидание ошибок равно нулю, $\mathbb{E}(u | X) = 0$;
4. Условная ковариационная матрица ошибок пропорциональна единичной, $\text{Var}(u | X) = \sigma^2 I$;
5. Оценка $\hat{\beta}$ получена методом наименьших квадратов, $\hat{\beta} = (X^T X)^{-1} X^T y$;

то

- (a) Оценка $\hat{\beta}$ является линейной по y ;
- (b) Оценка $\hat{\beta}$ является условно несмещённой, $\mathbb{E}(\hat{\beta} | X) = \beta$ и несмещённой, $\mathbb{E}(\hat{\beta}) = \beta$;
- (c) Оценка любого коэффициента $\hat{\beta}_j$ является наиболее эффективной в классе линейных несмещённых оценок.

Что означает «эффективная в классе линейных несмещённых оценок»? Это означает, что у любой другой линейной по y несмещённой оценки $\hat{\beta}_j^{\text{alt}}$ дисперсия не меньше, чем у МНК-оценки.

$$\text{Var}(\hat{\beta}_j | X) \leq \text{Var}(\hat{\beta}_j^{\text{alt}} | X).$$

Вывод теоремы можно усилить: для любой линейной комбинации коэффициентов $w^T \beta$ МНК-оценка $w^T \hat{\beta}$ эффективнее альтернативной оценки $w^T \hat{\beta}^{\text{alt}}$, то есть

$$\text{Var}(w^T \hat{\beta}_j | X) \leq \text{Var}(w^T \hat{\beta}_j^{\text{alt}} | X).$$

Доказательство. Линейность оценки по y видна прямо из её формулы, $\hat{\beta} = (X^T X)^{-1} X^T y$.

Проверим условную несмещённость,

$$\mathbb{E}(\hat{\beta} | X) = \mathbb{E}((X^T X)^{-1} X^T y | X) = (X^T X)^{-1} X^T \mathbb{E}(y | X).$$

Для удобства посчитаем $\mathbb{E}(y | X)$ отдельно,

$$\mathbb{E}(y | X) = \mathbb{E}(X\beta + u | X) = X\beta + \mathbb{E}(u | X) = X\beta.$$

И теперь завершаем вычисление $\mathbb{E}(\hat{\beta} | X)$:

$$\mathbb{E}(\hat{\beta} | X) = (X^T X)^{-1} X^T \mathbb{E}(y | X) = (X^T X)^{-1} X^T X \beta = \beta.$$

Мы доказали условную несмещённость оценки, $\mathbb{E}(\hat{\beta} \mid X) = \beta$. Безусловная несмещённость следует из свойства условного ожидания,

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(\mathbb{E}(\hat{\beta} \mid X)) = \mathbb{E}(\beta) = \beta.$$

Эффективность МНК-оценок — это реинкарнация теоремы Пифагора. Мы увидим, что дисперсия МНК-оценки — это квадрат длины катета, дисперсия альтернативной несмещённой оценки — квадрат длины гипотенузы.

Для примера рассмотрим оценку первого коэффициента бета, $\hat{\beta}_1$. Доказательство не меняется ни капли, если рассмотреть оценку другого коэффициента, скажем, $\hat{\beta}_7$ или даже оценку произвольной линейной комбинации коэффициентов бета, например, $\hat{\beta}_1 + \hat{\beta}_2 + \hat{\beta}_3$.

Итак, у нас есть две оценки, $\hat{\beta}_1$ и $\hat{\beta}_1^{\text{alt}}$. Обе они линейны по y , следовательно, $\hat{\beta}_1 = a^T y$ и $\hat{\beta}_1^{\text{alt}} = a_{\text{alt}}^T y$.

Замечаем, что $\text{Var}(\hat{\beta}_1 \mid X) = \sigma^2 a^T a$, и $\text{Var}(\hat{\beta}_1^{\text{alt}}) = \sigma^2 a_{\text{alt}}^T a_{\text{alt}}$. То есть дисперсии пропорциональны квадратам длин векторов a и a_{alt} . Осталось доказать, что вектор a не длиннее вектора a_{alt} :)

Для этого мы докажем, что вектор a_{alt} — это гипотенуза, а вектор a — катет. Нам нужно доказать, что вектор $a - a_{\text{alt}}$ перпендикулярен вектору a .

Разобьём доказательство перпендикулярности a и $a - a_{\text{alt}}$ на два шага:

Шаг 1. Вектор $a - a_{\text{alt}}$ перпендикулярен любому столбцу матрицы X .

Шаг 2. Вектор a является линейной комбинацией столбцов матрицы X .

здесь простая картинка с теоремой Пифагора!

Приступаем к шагу 1. Обе оценки несмещённые, поэтому для любых β должно выполняться:

$$\mathbb{E}(\hat{\beta}_1 \mid X) = \mathbb{E}(\hat{\beta}_1^{\text{alt}} \mid X)$$

Переносим всё в левую сторону:

$$\mathbb{E}((a^T - a_{\text{alt}}^T)(X\beta + u) \mid X) = 0$$

Получаем, что для любых β должно быть выполнено условие:

$$(a - a_{\text{alt}})^T X\beta = 0$$

Это возможно только, если вектор $(a - a_{\text{alt}})^T X$ равен нулю. Следовательно, вектор $(a - a_{\text{alt}})$ перпендикулярен любому столбцу X .

Приступаем к шагу 2.

Вспоминаем, что $\hat{\beta} = (X^T X)^{-1} X^T y = Ay$. Следовательно, нужная строка весов a^T — это первая строка в матрице $(X^T X)^{-1} X^T$. Замечаем, что выражение имеет вид $A \cdot X^T$.

Вспоминаем из линейной алгебры, что при умножении матриц AB получается матрица C , на которую можно взглянуть несколькими способами! Можно считать, что C — это разные линейные комбинации столбцов левой матрицы A . Можно считать, что C — это разные линейные комбинации строк правой матрицы B .

Применим второй взгляд :) Получаем, что строка a^T — линейная комбинация строк матрицы X^T . Или, другими словами, столбец a — линейная комбинация столбцов матрицы X . \square

Классическое доказательство эффективности, которое можно найти во многих учебниках, не замечает связи с теоремой Пифагора и исследует разницу ковариационных матриц. Приведём его здесь для демонстрации альтернативной техники!

Доказательство. У нас есть две линейных по y оценки: МНК-оценка и оценка-конкурент,

$$\hat{\beta} = (X^T X)^{-1} X^T y \text{ и } \hat{\beta}_{\text{alt}} = A_{\text{alt}}^T y.$$

Оценки ковариационных матриц этих оценок равны

$$\mathbb{V}\text{ar}(\hat{\beta} \mid X) = (X^T X)^{-1} \sigma^2 \text{ и } \mathbb{V}\text{ar}(\hat{\beta}_{\text{alt}} \mid X) = A^T A \sigma^2.$$

□

Условие несмещённости альтернативной оценки имеет вид

$$\mathbb{E}(\hat{\beta}_{\text{alt}} \mid X) = \mathbb{E}(A^T y \mid X) = A^T X \beta = \beta.$$

То есть для несмещённости альтернативной оценки должно выполняться условие $A^T X = I$. Для простоты рассмотрим случай $\sigma^2 = 1$. Мы докажем, что разность этих матриц $D = A^T A - (X^T X)^{-1}$ является положительно полуопределённой матрицы.

Вспомним из линейной алгебры определение и свойства положительно полуопределённой матрицы.

Определение 6.6 (положительно полуопределённая форма). Матрица D или квадратичная форма $q(v) = v^T D v$ называется положительно полуопределённой, если $q(v) \geq 0$ для любого вектора v .

Теорема 6.7 (свойства положительно полуопределённой матрицы). Матрица D является положительно полуопределённой, если и только если её можно записать в виде произведения $D = B^T B$.

У положительно полуопределённой матрицы D на диагонали находятся неотрицательные числа.

Если $D = A^T A - (X^T X)^{-1}$ — положительно полуопределена, то $d_{ii} \geq 0$ и, следовательно, $[A^T A]_{ii} \geq [(X^T X)^{-1}]_{ii}$, то есть, дисперсии альтернативных оценок не меньше дисперсий МНК-оценок.

Перейдём к доказательству положительной полуопределённости D :

Доказательство. Возьмём $B = A - X(X^T X)^{-1}$ и найдём $B^T B$:

$$B^T B = (A - X(X^T X)^{-1})^T (A - X(X^T X)^{-1}) = A^T A - A^T X(X^T X)^{-1} - (X^T X)^{-1} X^T A + (X^T X)^{-1} X^T X(X^T X)^{-1}$$

В силу несмещённости $A^T X = I$ или $X^T A = I$, поэтому

$$B^T B = A^T A - (X^T X)^{-1} - (X^T X)^{-1} + (X^T X)^{-1} = A^T A - (X^T X)^{-1}.$$

Мы видим, что матрица $D = A^T A - (X^T X)^{-1}$ оказалась разложенной в произведение $D = B^T B$ и, следовательно, матрица D положительно полуопределена. □

6.5. Статистические свойства остатков

Используя матричное представление для остатков $\hat{u} = My = Mu$, вычислим вектор математических ожиданий остатков

$$\mathbb{E}(\hat{u} \mid X) = \mathbb{E}(My \mid X) = M \mathbb{E}(y \mid X) = MX\beta = 0, \text{ так как } MX = 0.$$

(Ожидаемое значение остатков равно нулю, также как и ожидаемое значение ошибок, $\mathbb{E}(\hat{u} \mid X) = \mathbb{E}(u \mid X) = 0$.) и ковариационную матрицу остатков:

$$\mathbb{V}\text{ar}(\hat{u} \mid X) = \mathbb{V}\text{ar}(My \mid X) = M \mathbb{V}\text{ar}(y \mid X) M^T = M \sigma^2 I_n M^T = \sigma^2 M M^T = \sigma^2 M.$$

Вспомним, что у ковариационной матрицы ошибок $\text{Var}(u | X) = \sigma^2 I$ на диагонали стоят одинаковые элементы, а вне диагонали стоят нули. А у ковариационной матрицы остатков $\text{Var}(\hat{u} | X) = \sigma^2 M$ на диагоналях находятся разные элементы и вне диагонали элементы в общем случае не равны нулю.

Другими словами, остатки \hat{u}_i зависимы между собой и имеют разную дисперсию $\text{Var}(\hat{u}_i)$. Например, при наличии константы в регрессии остатки обязательно удовлетворяют соотношению $\sum \hat{u}_i = 0$.

Посчитаем ковариационную матрицу вектора остатков и вектора прогнозов:

$$\begin{aligned} \text{Cov}(\hat{u}, \hat{y} | X) &= \text{Cov}(Mu, Py | X) = \text{Cov}(Mu, P(X\beta + u) | X) = \text{Cov}(Mu, X\beta + Pu | X) = \\ &= \text{Cov}(Mu, Pu | X) = M \text{Cov}(u, u | X) P = M \sigma^2 I_n P^T = \sigma^2 MP = 0, \text{ так как } P^T = P \text{ и } MP = 0. \end{aligned}$$

Следовательно, вектор остатков и вектор прогнозов линейно независимы. Метод наименьших квадратов даёт наилучший линейный прогноз, то есть даже зная прогнозные значения \hat{y} нет возможности уменьшить остатки модели.

Посчитаем ковариационную матрицу вектора остатков и МНК-оценки вектора параметров β :

$$\begin{aligned} \text{Cov}(\hat{u}, \hat{\beta} | X) &= \text{Cov}(Mu, \beta + (X^T X)^{-1} X^T u | X) = \text{Cov}(Mu, (X^T X)^{-1} X^T u | X) = \\ &= M \text{Cov}(u, u | X) X (X^T X)^{-1} = M \text{Cov}(u, u | X) X (X^T X)^{-1} = \sigma^2 M X (X^T X)^{-1}, \text{ так как } MX = 0. \end{aligned}$$

Следовательно, вектор остатков и вектор МНК-оценок параметров модели.

6.6. Оценивание дисперсии

Метод наименьших квадратов позволяет оценить вектор параметров β , однако ~~власти скрывают настоящую дисперсию~~ никак не оценивает неизвестный параметр σ^2 . Интуиция говорит, что высокая дисперсия ошибок u_i должна проявляться в высоком разбросе \hat{u}_i , поэтому разумно попробовать построить оценку $\hat{\sigma}^2$ на базе $RSS = \sum \hat{u}_i^2$.

Для построения оценки $\hat{\sigma}^2$ найдём ожидание $\mathbb{E}(RSS | X)$:

Теорема 6.8 (ожидание суммы квадратов остатков). Если выполнены предпосылки теоремы Гаусса — Маркова,

1. Модель линейна по параметрам: $y = X\beta + u$;
2. Матрица X размера $[n \times k]$ имеет полный ранг k .
3. Условное ожидание ошибок равно нулю, $\mathbb{E}(u | X) = 0$;
4. Условная ковариационная матрица ошибок пропорциональна единичной, $\text{Var}(u | X) = \sigma^2 I$;
5. Оценка $\hat{\beta}$ получена методом наименьших квадратов, $\hat{\beta} = (X^T X)^{-1} X^T y$;

то $\mathbb{E}(RSS | X) = \mathbb{E}(\sum \hat{u}_i^2 | X) = (n - k)\sigma^2$.

Из этой теоремы следует, что оценка $\hat{\sigma}^2 = RSS/(n - k)$ — несмещённая оценка для неизвестной дисперсии σ^2 .

Доказательство. На помощь нам придёт след матрицы! След матрицы прекрасен двумя свойствами. Во-первых, его можно менять местами с математическим ожиданием, $\mathbb{E}(\text{trace } W) = \text{trace } \mathbb{E}(W)$. Во-вторых, внутри следа можно переставлять местами перемножаемые матрицы, $\text{trace}(AB) = \text{trace}(BA)$.

Кроме того, на скалярную величину след можно навесить совершенно бесплатно! Если величина R — не вектор, а скаляр, то $\text{trace } R = R$.

Продолжаем,

$$\mathbb{E}(\hat{u}^T \hat{u} \mid X) = \mathbb{E}(\text{trace}(\hat{u}^T \hat{u}) \mid X) = \mathbb{E}(\text{trace}(\hat{u} \hat{u}^T) \mid X) = \text{trace } \mathbb{E}(\hat{u} \hat{u}^T \mid X).$$

Подумаем о середине,

$$\mathbb{E}(\hat{u} \hat{u}^T \mid X) = \mathbb{E}(Mu(Mu)^T \mid X) = \mathbb{E}(Mu u^T M^T \mid X) = M \mathbb{E}(u u^T \mid X) M^T.$$

Вспомним, что матрица M — проектор, поэтому $M^T = M$, $M^2 = M$. У матрицы $u u^T$ на диагонали стоят u_i^2 , вне диагонали — $u_i u_j$. Поэтому $\mathbb{E}(u u^T \mid X) = \sigma^2 I$. Завершаем вычисления,

$$\mathbb{E}(\hat{u} \hat{u}^T \mid X) = M \mathbb{E}(u u^T \mid X) M^T = M \cdot \sigma^2 I \cdot M^T = \sigma^2 M^2 = \sigma^2 M$$

След проектора равен размерности пространства, на которое он проецирует, поэтому $\text{trace } M = n - k$ и

$$\mathbb{E}(RSS \mid X) = \text{trace}(\sigma^2 M) = (n - k) \sigma^2$$

И мы легко строим несмещённую оценку, $\hat{\sigma}^2 = RSS / (n - k)$,

$$\mathbb{E}(\hat{\sigma}^2 \mid X) = \mathbb{E}\left(\frac{RSS}{n - k} \mid X\right) = \frac{(n - k) \sigma^2}{n - k} = \sigma^2$$

□

Выборочная дисперсия при случайной выборке

Заметим, что данная теорема обобщает старый факт про выборочную дисперсию! Вспомним, что для выборки из независимых y_i с ожиданием $\mathbb{E}(y_i) = \mu$ и дисперсией $\text{Var}(y_i) = \sigma^2$ несмещённая оценка дисперсии имеет вид

$$\hat{\sigma}^2 = \frac{\sum (y_i - \bar{y})^2}{n - 1}.$$

В данном случае величины y_i можно представить в виде $y_i = \mu + u_i$. Тогда предпосылки теоремы Гаусса — Маркова выполнены, матрица регрессоров X — это просто единственный столбец-регрессор из единиц, $k = 1$, $\beta = \mu$. В этом случае $\hat{\beta} = \bar{y}$, все прогнозы равны $\hat{u}_i = \bar{y}$ и $RSS = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \bar{y})^2$. И мы видим, что новая оценка совпадает в этом случае со старой:

$$\hat{\sigma}^2 = \frac{RSS}{n - k} = \frac{\sum (y_i - \hat{y}_i)^2}{n - 1} = \frac{\sum (y_i - \bar{y})^2}{n - 1}$$

Оценка дисперсии оценок коэффициентов

Для построения доверительных интервалов для коэффициентов β_j нам понадобятся оценки дисперсий $\text{Var}(\hat{\beta}_j \mid X)$. К счастью, у нас есть несмещённая оценка $\hat{\sigma}^2$ для σ^2 . Из неё мы легко построим оценку и для неизвестной ковариационной матрицы $\text{Var}(\hat{\beta} \mid X) = \sigma^2 (X^T X)^{-1}$. А именно, мы просто подставим оценку дисперсии вместо неизвестной дисперсии:

$$\widehat{\text{Var}}(\hat{\beta} \mid X) = \hat{\sigma}^2 (X^T X)^{-1} = \frac{RSS}{n - k} (X^T X)^{-1}.$$

Уточним, что эту оценку мы вывели из предпосылок теоремы Гаусса — Маркова. Если использовать другие предпосылки, то ковариационная матрица $\text{Var}(\hat{\beta} \mid X)$ перестанет быть равной $\sigma^2(X^T X)^{-1}$ и нам потребуется другой способ оценивания.

Оценки корней из дисперсий оценок называются стандартными ошибками оценок (standard errors).

$$\text{se}(\hat{\beta}_j) = \sqrt{\widehat{\text{Var}}(\hat{\beta}_j \mid X)} = \sqrt{\hat{\sigma}^2[(X^T X)^{-1}]_{jj}}$$

Число $[(X^T X)^{-1}]_{jj}$ мы берём в матрице $(X^T X)^{-1}$ из j -й строки j -го столбца.

6.7. Неправильная спецификация модели

Одной из предпосылок теоремы Гаусса — Маркова является правильный выбор спецификации, при котором мы регрессируем y в точности на набор истинных регрессоров. В реальности такое условие вряд ли выполнимо, так как не до всех регрессоров мы способны догадаться. А если догадаемся, то не все сможем измерить или собрать. Можно ли допустить неполную спецификацию модели, но получить BLUE-оценки (несмещённые и эффективные в классе линейных) для собранных регрессоров?

Рассмотрим для начала случай, когда при оценивании модели мы пропускаем часть важных регрессоров. Истинная модель имеет вид

$$y = W\beta + V\gamma + u,$$

где W — матрица регрессоров размерности $[n \times k_1]$, V — матрица регрессоров размерности $[n \times k_2]$. Обозначим через $X = [W \ V]$ $[n \times k]$ матрицу всех регрессоров, где $k = k_1 + k_2$.

Вместо истинной модели оценивается следующая модель:

$$y = W\beta + \nu,$$

где ν — вектор случайных ошибок в оцениваемой модели.

Утверждение 6.1. тут должно быть утверждение про смещённость

Доказательство. Пусть X — истинный набор регрессоров, а W — собранный датасет. При этом $X = [W \ V]$. Тогда новая МНК-оценка получается из изменившейся предпосылки о правильности спецификации $\tilde{\beta} = (W^T W)^{-1} W^T y$.

Предположим истинную модель можно записать как

$$y = W\beta_W + V\beta_V + u,$$

Мы бы всё равно хотели получать несмещённую оценку.

$$\begin{aligned} \mathbb{E}(\tilde{\beta} \mid W) &= (W^T W)^{-1} W^T \mathbb{E}(y \mid W) = (W^T W)^{-1} W^T \mathbb{E}(W\beta_W + V\beta_V + u \mid W) = (W^T W)^{-1} W^T (W\beta_W + \mathbb{E}(V\beta_V \mid W)) \\ &= (W^T W)^{-1} W^T (W\beta_W + V\beta_V) = \beta_W + (W^T W)^{-1} W^T V \beta_V. \end{aligned}$$

Отсюда видно, что в общем случае $\mathbb{E}(\tilde{\beta} \mid W) \neq \beta_W$: возникает смещение. Оценка не смещенна только тогда, когда $W^T V = 0$ (то есть включённые и пропущенные регрессоры ортогональны) либо $\beta_V = 0$ (пропущенные переменные вообще не нужны и не влияют на y). \square

Утверждение 6.2. Пусть $\text{Var}(\hat{\beta} \mid W, V)$ — ковариационная матрица вектора оценок $\hat{\beta}$, полученного по полному набору регрессоров $X = [W \ V]$, а $\text{Var}(\tilde{\beta} \mid W)$ — ковариационная матрица вектора оценок $\tilde{\beta}$, полученного по регрессорам из матрицы W . Тогда матрица $\text{Var}(\hat{\beta} \mid W, V) - \text{Var}(\tilde{\beta} \mid W)$ является положительно полуопределённой матрицей.

Утверждение 6.2 означает, что на диагонали матрицы $\text{Var}(\hat{\beta} \mid W, V) - \text{Var}(\tilde{\beta} \mid W)$ стоят неотрицательные значения. В свою очередь, диагональный элемент с индексами jj представляет собой разницу дисперсий оценок коэффициента β_j , полученных по полному и по сокращенному набору переменных. Это означает, что $\text{Var}(\hat{\beta}_j \mid W, V) - \text{Var}(\tilde{\beta}_j \mid V) \geq 0$, то есть оценка $\tilde{\beta}_j$ имеет меньшую условную дисперсию. Из-за меньшей условной дисперсии оценка $\tilde{\beta}_j$ может получиться более эффективной по сравнению с оценкой $\hat{\beta}_j$.

Утверждение 6.3. Оценка дисперсии случайной ошибки $\tilde{\sigma}^2 = \frac{RSS}{n-k_1}$, полученная по модели с пропущенными переменными, является смещённой,

$$\mathbb{E}(\tilde{\sigma}^2 \mid W) \neq \sigma^2.$$

Доказательство. Вспомним матричное представление RSS :

$$RSS = y^T M y, \text{ где } M = I_n - W(W^T W)^{-1} W^T.$$

Рассчитаем математическое ожидание RSS , учитывая, что истинной моделью является модель по набору регрессоров $X = [W \ V]$:

$$\begin{aligned} \mathbb{E}(RSS \mid W, V) &= \mathbb{E}(y^T M y \mid W, V) = \mathbb{E}((W\beta + V\gamma + u)^T M (W\beta + V\gamma + u) \mid W, V) = \\ &= \mathbb{E}(u^T M u + 2\gamma^T V^T M u + \gamma^T V^T M V \gamma \mid W, V) = \sigma^2(n - k_1) + \gamma^T V^T M V \gamma. \end{aligned}$$

Выше мы использовали следующие результаты:

- $MW = 0$;
- $\mathbb{E}(2\gamma^T V^T M u \mid W, V) = 2\gamma^T V^T M \mathbb{E}(u \mid W, V) = 0$;
- $\mathbb{E}(u^T M u \mid W, V) = \sigma^2(n - k_1)$.

Таким образом, получаем, что

$$\mathbb{E}(\tilde{\sigma}^2 \mid W, V) = \mathbb{E}\left(\frac{RSS}{n - k_1} \mid W, V\right) = \frac{1}{n - k_1}(\sigma^2(n - k_1) + \gamma^T V^T M V \gamma) = \sigma^2 + \frac{1}{n - k_1} \gamma^T V^T M V \gamma.$$

Оценка дисперсии $\tilde{\sigma}^2$ будет несмещённой только, если $\gamma = 0$. Равенство $\gamma = 0$ означает, что пропущенных переменных нет и $X = W$. Заметим также, что $V^T M V = (M V)^T M V$, что означает, что матрица $V^T M V$ является положительно полуопределённой. Следовательно, смещение оценки $\tilde{\sigma}^2$ в общем случае положительное. \square

Далее проанализируем, что происходит со свойствами несмещённости и эффективности МНК-оценок при включении в модель лишних регрессоров.

Теперь истинной моделью является

$$y = X\beta + u.$$

Вместо истинной модели оценивается следующая модель:

$$y = X\beta + R\gamma + \nu,$$

где R — матрица лишних регрессоров.

Утверждение 6.4. При включении лишних регрессоров МНК-оценка $\tilde{\beta}$, полученная в модели с набором регрессоров $(X \ R)$, остаётся несмещённой, то есть $\mathbb{E}(\tilde{\beta} \mid X, R) = \beta$.

Утверждение 6.5. Пусть $\text{Var}(\hat{\beta} \mid X)$ — ковариационная матрица вектора оценок $\hat{\beta}$, полученного по истинному набору регрессоров X , а $\text{Var}(\tilde{\beta} \mid X, R)$ — ковариационная матрица вектора оценок $\tilde{\beta}$, полученного по регрессорам из матрицы

$$(X \ R)$$

Тогда матрица $\text{Var}(\tilde{\beta} \mid X, R) - \text{Var}(\hat{\beta} \mid X)$ является положительно полуопределённой матрицей.

Утверждение 6.5 означает, что на диагонали матрицы $\text{Var}(\tilde{\beta} \mid X, R) - \text{Var}(\hat{\beta} \mid X)$ стоят неотрицательные значения. В свою очередь, диагональный элемент с индексами jj представляет собой разницу дисперсий оценок коэффициента β_j , полученных по расширенному и по истинному наборам переменных. Это означает, что $\text{Var}(\tilde{\beta}_j \mid X, R) - \text{Var}(\hat{\beta}_j \mid X) \geq 0$, то есть оценка $\tilde{\beta}_j$ имеет большую условную дисперсию. Из-за большей условной дисперсии оценка $\tilde{\beta}_j$ может получиться менее эффективной по сравнению с оценкой $\hat{\beta}_j$.

Утверждение 6.6. Оценка дисперсии случайной ошибки $\tilde{\sigma}^2 = \frac{RSS}{n-k}$, полученная по модели с лишними регрессорами, является несмещённой,

$$\mathbb{E}(\tilde{\sigma}^2 \mid W) = \sigma^2.$$

6.8. Задачи для семинара:

Задача 23. Исследовательница Мишель собрала данные по 20 студентам. Переменная y_i — количество решённых задач по эконометрике i -м студентом, а x_i — количество просмотренных серий любимого сериала за прошедший год. Оказалось, что $\sum y_i = 10$, $\sum x_i = 0$, $\sum x_i^2 = 40$, $\sum y_i^2 = 50$, $\sum x_i y_i = 60$.

- Найдите МНК-оценки коэффициентов парной регрессии с константой.
- В рамках предположения $\mathbb{E}(u_i \mid X) = 0$ найдите $\mathbb{E}(y_i \mid X)$, $\mathbb{E}(\hat{\beta}_j \mid X)$, $\mathbb{E}(\hat{u}_i \mid X)$, $\mathbb{E}(\hat{y}_i \mid X)$.
- Предположим дополнительно, что $\text{Var}(u_i \mid X) = \sigma^2$ и u_i при фиксированных X независимы. Найдите $\text{Var}(y_i \mid X)$, $\text{Var}(y_i(x_i - \bar{x}) \mid X)$, $\text{Var}(\sum y_i(x_i - \bar{x}) \mid X)$, $\text{Var}(\hat{\beta}_2 \mid X)$.

Решение. а) Модель парной регрессии с константой:

$$y_i = \beta_1 + \beta_2 x_i + u_i.$$

МНК-оценки коэффициентов имеют вид:

$$\hat{\beta}_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}, \quad \hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}.$$

Вычислим следующие величины:

$$\bar{x} = \frac{\sum x_i}{n} = \frac{0}{20} = 0, \quad \bar{y} = \frac{\sum y_i}{n} = \frac{10}{20} = 0.5,$$

$$\sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - \bar{y} \sum x_i = 60 - 0.5 \cdot 0 = 60,$$

$$\sum (x_i - \bar{x})^2 = \sum x_i^2 - n\bar{x}^2 = 40 - 20 \cdot 0 = 40.$$

Проделаем финальные вычисления:

$$\hat{\beta}_2 = \frac{60}{40} = 1.5, \quad \hat{\beta}_1 = 0.5 - 1.5 \cdot 0 = 0.5.$$

б) При $\mathbb{E}(u_i | X) = 0$:

- (a) $\mathbb{E}(y_i | X) = \beta_1 + \beta_2 x_i$
- (b) $\mathbb{E}(\hat{\beta}_j | X) = \beta_j$ для $j = 1, 2$
- (c) $\mathbb{E}(\hat{u}_i | X) = 0$
- (d) $\mathbb{E}(\hat{y}_i | X) = \beta_1 + \beta_2 x_i$

в) При $\text{Var}(u_i | X) = \sigma^2$ и независимости u_i :

- (a) $\text{Var}(y_i | X) = \sigma^2$
- (b) $\text{Var}(y_i(x_i - \bar{x}) | X) = \sigma^2 x_i^2$ (т.к. $\bar{x} = 0$)
- (c) $\text{Var}(\sum y_i(x_i - \bar{x}) | X) = 40\sigma^2$
- (d) $\text{Var}(\hat{\beta}_2 | X) = \frac{\sigma^2}{40}$

Задача 24. Рассмотрим модель $y_i = \beta_1 + \beta_2 x_i + \beta_3 z_i + u_i$, где

$$X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}, \quad y = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}, \quad u = \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \end{pmatrix}.$$

Случайные ошибки u_i независимы и нормально распределены с $\mathbb{E}(u | X) = 0$ и $\text{Var}(u | X) = \sigma^2 I$.

Для удобства расчётов даны матрицы: $X^T X$, $(X^T X)^{-1}$ и $X^T y$:

$$X^T X = \begin{pmatrix} 5 & 3 & 1 \\ 3 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix}, \quad (X^T X)^{-1} = \begin{pmatrix} 0.5 & -0.5 & 0 \\ -0.5 & 1 & -0.5 \\ 0 & -0.5 & 1.5 \end{pmatrix}, \quad X^T y = \begin{pmatrix} 15 \\ 12 \\ 5 \end{pmatrix}.$$

- а) Определите n и k .
- б) Вычислите МНК оценку вектора β .
- в) Найдите $\hat{\sigma}^2$, $\mathbb{E}(\hat{\sigma}^2 | X)$.
- г) Найдите $\text{Var}(u_1)$, $\text{Var}(\beta_1)$, $\text{Var}(\hat{\beta}_1 | X)$, $\widehat{\text{Var}}(\hat{\beta}_1 | X)$, $\mathbb{E}(\hat{\beta}_1^2 | X) - \beta_1^2$;
- д) Найдите $\text{Cov}(\hat{\beta}_2, \hat{\beta}_3 | X)$, $\widehat{\text{Cov}}(\hat{\beta}_2, \hat{\beta}_3 | X)$, $\text{Var}(\hat{\beta}_2 - \hat{\beta}_3 | X)$, $\widehat{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3 | X)$;
- е) Найдите $\text{Var}(\beta_2 - \beta_3)$, $\text{Corr}(\hat{\beta}_2, \hat{\beta}_3 | X)$, $\widehat{\text{Corr}}(\hat{\beta}_2, \hat{\beta}_3 | X)$;

Решение. а) Число наблюдений $n = 5$. Число регрессоров, включая свободный член равно $k = 3$.

б) МНК-оценка вектора β равна $\hat{\beta} = (X^T X)^{-1} X^T y$. Тогда

$$\hat{\beta} = \begin{pmatrix} 0.5 & -0.5 & 0 \\ -0.5 & 1 & -0.5 \\ 0 & -0.5 & 1.5 \end{pmatrix} \times \begin{pmatrix} 15 \\ 12 \\ 5 \end{pmatrix} = \begin{pmatrix} 1.5 \\ 2 \\ 1.5 \end{pmatrix}$$

в) Несмещенная оценка дисперсии случайной ошибки равна $\hat{\sigma}^2 = \frac{RSS}{n-k} = \frac{RSS}{5-3}$. Вычислим RSS . Знаем, что $RSS = y^T (I - X(X^T X)^{-1} X^T) y = 1$.

Тогда $\hat{\sigma}^2 = 1/2$.

Так как по построению оценка $\hat{\sigma}^2$ несмещённая, то $\mathbb{E}(\hat{\sigma}^2 | X) = \sigma^2$.

г)

$$\mathbb{V}\text{ar}(u_1) = \mathbb{V}\text{ar}(u)_{(1,1)} = \sigma^2 \cdot I_{(1,1)} = \sigma^2$$

$$\mathbb{V}\text{ar}(\beta_1) = 0,$$

$$\mathbb{V}\text{ar}(\hat{\beta}_1) = \sigma^2 (X'X)^{-1}_{(1,1)} = 0.5\sigma^2$$

$$\widehat{\mathbb{V}\text{ar}}(\hat{\beta}_1) = \hat{\sigma}^2 (X'X)^{-1}_{(1,1)} = 0.5\hat{\sigma}^2_{(1,1)} = 0.5 \frac{RSS}{5-3} = 0.25RSS = 0.25y'(I - X(X'X)^{-1}X')y = 0.25 \cdot 1 = 0.25$$

Так как оценки МНК являются несмещёнными, то $\mathbb{E}(\hat{\beta}) = \beta$, значит:

$$\mathbb{E}(\hat{\beta}_1) - \beta_1^2 = \mathbb{E}(\hat{\beta}_1) - (\mathbb{E}(\hat{\beta}_1))^2 = \widehat{\mathbb{V}\text{ar}}(\hat{\beta}_1) = 0.25$$

д)

$$\mathbb{C}\text{ov}(\hat{\beta}_2, \hat{\beta}_3) = \sigma^2 (X'X)^{-1}_{(2,3)} = \sigma^2 \cdot \left(-\frac{1}{2}\right)$$

$$\widehat{\mathbb{C}\text{ov}}(\hat{\beta}_2, \hat{\beta}_3) = \widehat{\mathbb{V}\text{ar}}(\hat{\beta})_{(2,3)} = \hat{\sigma}^2 (X'X)^{-1}_{(2,3)} = \frac{1}{2} \cdot \left(-\frac{1}{2}\right) = -\frac{1}{4}$$

$$\begin{aligned} \mathbb{V}\text{ar}(\hat{\beta}_2 - \hat{\beta}_3) &= \mathbb{V}\text{ar}(\hat{\beta}_2) + \mathbb{V}\text{ar}(\hat{\beta}_3) + 2\mathbb{C}\text{ov}(\hat{\beta}_2, \hat{\beta}_3) = \\ &= \sigma^2((X'X)^{-1}_{(2,2)} + (X'X)^{-1}_{(3,3)} + 2(X'X)^{-1}_{(2,3)}) = \sigma^2(1 + 1.5 + 2 \cdot (-0.5)) = 1.5\sigma^2 \end{aligned}$$

$$\begin{aligned} \widehat{\mathbb{V}\text{ar}}(\hat{\beta}_2 - \hat{\beta}_3) &= \widehat{\mathbb{V}\text{ar}}(\hat{\beta}_2) + \widehat{\mathbb{V}\text{ar}}(\hat{\beta}_3) + 2\widehat{\mathbb{C}\text{ov}}(\hat{\beta}_2, \hat{\beta}_3) = \\ &= \hat{\sigma}^2((X'X)^{-1}_{(2,2)} + (X'X)^{-1}_{(3,3)} + 2(X'X)^{-1}_{(2,3)}) = \frac{1}{2} \cdot 1.5 = 0.75 \end{aligned}$$

е)

$$\mathbb{V}\text{ar}(\beta_2 - \beta_3) = 0$$

$$\text{Corr}(\hat{\beta}_2, \hat{\beta}_3) = \frac{\mathbb{C}\text{ov}(\hat{\beta}_2, \hat{\beta}_3)}{\sqrt{\mathbb{V}\text{ar}(\hat{\beta}_2) \mathbb{V}\text{ar}(\hat{\beta}_3)}} = \frac{-0.5}{\sqrt{1 \cdot 1.5}} = -\frac{1}{\sqrt{6}}$$

$$\widehat{\text{Corr}}(\hat{\beta}_2, \hat{\beta}_3) = \frac{\widehat{\mathbb{C}\text{ov}}(\hat{\beta}_2, \hat{\beta}_3)}{\sqrt{\widehat{\mathbb{V}\text{ar}}(\hat{\beta}_2) \widehat{\mathbb{V}\text{ar}}(\hat{\beta}_3)}} = \frac{-0.5}{\sqrt{1 \cdot 1.5}} = -\frac{1}{\sqrt{6}}$$

Задача 25. Рассмотрим классическую линейную модель $y = X\beta + u$ с предпосылками Гаусса — Маркова: $\mathbb{E}(u | X) = 0$ и $\text{Var}(u | X) = \sigma^2 I$. Для всех случайных векторов $(y, \hat{\beta}, \hat{y}, u, \hat{u}, \bar{y})$ найдите все возможные условные математические ожидания и ковариационные матрицы. $\mathbb{E}(\cdot)$, $\text{Var}(\cdot)$, $\text{Cov}(\cdot, \cdot)$.

Решение. Для начала вычислим вектора математических ожиданий:

- а) $\mathbb{E}(y | X) = \mathbb{E}(X\beta + u | X) = X\beta + \mathbb{E}(u | X) = X\beta$.
- б) $\mathbb{E}(\hat{\beta} | X) = \mathbb{E}((X^T X)^{-1} X^T y | X) = \mathbb{E}((X^T X)^{-1} X^T (X\beta + u) | X) = \beta + (X^T X)^{-1} X^T \mathbb{E}(u | X) = \beta$.
- в) $\mathbb{E}(\hat{y} | X) = \mathbb{E}(X\hat{\beta} | X) = X \mathbb{E}(\hat{\beta} | X) = X\beta$.
- г) $\mathbb{E}(u | X) = 0$.
- д) $\mathbb{E}(\hat{u} | X) = \mathbb{E}(Mu | X) = M \mathbb{E}(u | X) = 0$.

Теперь рассчитаем для указанных векторов ковариационные матрицы:

- $\text{Var}(y | X) = \text{Var}(X\beta + u | X) = \text{Var}(u | X) = \sigma^2 I_n$
- $\text{Var}(\hat{\beta} | X) = \text{Var}((X^T X)^{-1} X^T y | X) = \text{Var}((X^T X)^{-1} X^T (X\beta + u) | X) = \text{Var}(\beta + (X^T X)^{-1} X^T u | X) = \text{Var}((X^T X)^{-1} X^T u | X) = (X^T X)^{-1} X^T \text{Var}(u) X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$
- $\text{Var}(\hat{y} | X) = \text{Var}(Hy | X) = H \text{Var}(y | X) H^T = H \sigma^2 I_n H^T = \sigma^2 H H^T = \sigma^2 H = \sigma^2 X (X^T X)^{-1} X^T$
- $\text{Var}(u | X) = \sigma^2 I$
- $\text{Var}(\hat{u} | X) = \text{Var}(Mu | X) = M \text{Var}(u | X) M^T = M \sigma^2 I M^T = \sigma^2 M M^T = \sigma^2 (I - X (X^T X)^{-1} X^T)$
- $\text{Var}(\bar{y} | X) = \text{Var}(\pi y | X) = \pi \text{Var}(y | X) \pi^T = \sigma^2 \pi \pi^T = \sigma^2 \cdot s(s^T s)^{-1} s^T$
- $\text{Cov}(y, \hat{\beta} | X) = \text{Cov}(y, (X^T X)^{-1} X^T y | X) = (X^T X)^{-1} X^T \text{Var}(y) = \sigma^2 (X^T X)^{-1} X^T$
- $\text{Cov}(y, \hat{y} | X) = \text{Cov}(y, Hy | X) = \text{Var}(y | X) H^T = \sigma^2 X (X^T X)^{-1} X^T$
- $\text{Cov}(y, u | X) = \text{Cov}(X\beta + u, u | X) = \text{Cov}(u, u | X) = \text{Var}(u | X) = \sigma^2 I_n$
- $\text{Cov}(y, \hat{u} | X) = \text{Cov}(y, My | X) = \text{Var}(y | X) M^T = \sigma^2 M = \sigma^2 (I - X (X^T X)^{-1} X^T)$
- $\text{Cov}(y, \bar{y} | X) = \text{Cov}(y, \pi y | X) = \text{Cov}(y, y | X) \pi^T = \sigma^2 (s(s^T s)^{-1} s^T)$
- $\text{Cov}(\hat{\beta}, \hat{y} | X) = \text{Cov}((X^T X)^{-1} X^T y, Hy | X) = (X^T X)^{-1} X^T \text{Var}(y) H^T = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} X^T = \sigma^2 (X^T X)^{-1} X^T$
- $\text{Cov}(\hat{\beta}, u | X) = \text{Cov}((X^T X)^{-1} X^T y, u | X) = \text{Cov}((X^T X)^{-1} X^T (X\beta + u), u | X) = \text{Var}(u | X) = \sigma^2 I_n$
- $\text{Cov}(\hat{\beta}, \hat{u} | X) = \text{Cov}((X^T X)^{-1} X^T y, Hu | X) = (X^T X)^{-1} X^T \text{Cov}(y, u | X) H^T = (X^T X)^{-1} X^T \sigma^2 H = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} X^T = \sigma^2 (X^T X)^{-1} X^T$
- $\text{Cov}(\hat{\beta}, \bar{y} | X) = \text{Cov}(\hat{\beta}, \pi y | X) = \text{Cov}((X^T X)^{-1} X^T y, \pi y | X) = (X^T X)^{-1} X^T \text{Var}(y | X) \pi^T = \sigma^2 (X^T X)^{-1} X^T (s(s^T s)^{-1} s^T)$
- $\text{Cov}(\hat{y}, u | X) = \text{Cov}(Hy, u | X) = \text{Cov}(Hu, u | X) = H \text{Cov}(u, u | X) = H \text{Var}(u | X) = \sigma^2 H = \sigma^2 X (X^T X)^{-1} X^T$

- $\text{Cov}(\hat{y}, \hat{u} \mid X) = \text{Cov}(Hy, My \mid X) = H \text{Var}(y \mid X) M^T = \sigma^2 H M = 0$
- $\text{Cov}(\hat{y}, \bar{y} \mid X) = \text{Cov}(Hy, \pi y) = H \text{Var}(y \mid X) \pi^T = \sigma^2 (X^T X)^{-1} X^T (s(s^T s)^{-1} s^T)$
- $\text{Cov}(u, \hat{u} \mid X) = \text{Cov}(u, Mu \mid X) = \text{Var}(u \mid X) M^T = \sigma^2 M = \sigma^2 (I - X(X^T X)^{-1} X^T)$
- $\text{Cov}(u, \bar{y} \mid X) = \text{Cov}(u, \pi y \mid X) = \text{Cov}(u, \pi(X\beta + u) \mid X) = \text{Cov}(u, \pi u \mid X) = \text{Var}(u \mid X) \pi = \sigma^2 \pi = \sigma^2 (s(s^T s)^{-1} s^T)$
- $\text{Cov}(\hat{u}, \bar{y} \mid X) = \text{Cov}(My, \pi y \mid X) = M \text{Var}(u \mid X) \pi^T = \sigma^2 M \pi = 0$

Задача 26. Рассмотрим модель $y_i = \beta x_i + u_i$ с двумя наблюдениями, $x_1 = 1, x_2 = 2$. Величины u_1 и u_2 независимы и равновероятно равны $+1$ или -1 .

- Найдите оценку $\hat{\beta}_{\text{ols}}$ для β с помощью метода наименьших квадратов.
- Чему равна дисперсия $\text{Var}(\hat{\beta}_{\text{ols}} \mid x)$ и ожидание $\mathbb{E}(\hat{\beta}_{\text{ols}} \mid x)$?
- Постройте несмещённую оценку $\hat{\beta}_{\text{best}}$ с наименьшей дисперсией.
- Чему равна дисперсия $\text{Var}(\hat{\beta}_{\text{best}} \mid x)$?
- А как же теорема Гаусса — Маркова? Почему в данном примере удаётся построить оценку с дисперсией меньше, чем у оценки методом наименьших квадратов?

Решение. а) $\hat{\beta}_{\text{ols}} = (y_1 + 2y_2)/5$;

б) $\text{Var}(\hat{\beta}_{\text{ols}} \mid x) = 1/5$;

в) Заметим, что по величине $2y_1 - y_2$ можно однозначно восстановить величины ошибок u_1 и u_2 . Например, если $2y_1 - y_2 = 3$, то $u_1 = 1, u_2 = -1$.

$$\hat{\beta}_{\text{best}} = \begin{cases} y_1 + 1, & \text{если } 2y_1 - y_2 < 0, \\ y_1 - 1, & \text{если } 2y_1 - y_2 > 0. \end{cases}$$

г) Шок контент, $\text{Var}(\hat{\beta}_{\text{best}} \mid x) = 0$.

д) Построенная оценка $\hat{\beta}_{\text{best}}$ является нелинейной по y , а теорема Гаусса — Маркова гарантирует только, что метод наименьших квадратов порождает несмещённую оценку с наименьшей дисперсией среди линейных по y оценок.

Задача 27. Пусть r — случайный вектор размерности $[n \times 1]$, а A — неслучайная матрица подходящей размерности. Докажите, что справедлива следующая формула для математического ожидания квадратичной формы:

$$\mathbb{E}(r^T A r) = \text{trace}(A \text{Var}(r)) + \mathbb{E}(r^T) A \mathbb{E}(r).$$

Решение. Распишем математическое ожидание, используя свойства следа:

$$\mathbb{E}(r^T A r) = \mathbb{E}(\text{trace}(r^T A r)) = \mathbb{E}(\text{trace}(A r r^T)) = \mathbb{E}(\text{trace}(A r r^T)) = \text{trace} \mathbb{E}(A r r^T) = \text{trace}(A \mathbb{E}(r r^T)).$$

Вспомним, что

$$\text{Var}(r) = \mathbb{E}(r r^T) - \mathbb{E}(r) \mathbb{E}(r^T),$$

откуда получаем:

$$\begin{aligned}\mathbb{E}(r^T A r) &= \text{trace}(A \mathbb{E}(r r^T)) = \text{trace}(A(\text{Var}(r) + \mathbb{E}(r) \mathbb{E}(r^T))) = \\ &= \text{trace}(A \text{Var}(r)) + \text{trace}(A \mathbb{E}(r) \mathbb{E}(r^T)) = \text{trace}(A \text{Var}(r)) + \text{trace}(\mathbb{E}(r^T) A \mathbb{E}(r)) = \\ &= \text{trace}(A \text{Var}(r)) + \mathbb{E}(r^T) A \mathbb{E}(r).\end{aligned}$$

Задача 28. Предположим, что все предпосылки теоремы Гаусса – Маркова выполнены. Вычислите математические ожидания для TSS , ESS и RSS , используя их матричные представления.

Решение. Вспомним матричное представление для TSS , ESS и RSS :

$$TSS = y^T (I - \pi) y,$$

$$ESS = y^T (H - \pi) y,$$

$$RSS = y^T (I - H) y.$$

Заметим, что все три показателя представлены в виде квадратичной формы. В терминах предыдущей задачи $r = y$, а матрица A равна $(I - \pi)$ для TSS , $(H - \pi)$ для ESS и $(I - H)$ для RSS . Здесь матрица I имеет размерность $[n \times n]$.

Предварительно напомним, что

$$\text{trace}(I) = n, \quad \text{trace}(H) = k.$$

Из определения матрицы π

$$\pi = \frac{1}{n} \begin{pmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \cdots & \cdots & \cdots & \cdots \\ 1 & 1 & \cdots & 1 \end{pmatrix}$$

нетрудно убедиться, что $\text{trace}(\pi) = 1$.

Используя формулу для математического ожидания квадратичной формы, полученную в предыдущем упражнении, можно легко посчитать интересующие нас математические ожидания.

Начнём с $\mathbb{E}(TSS)$:

$$\begin{aligned}\mathbb{E}(TSS) &= \mathbb{E}(y^T (I - \pi) y) = \text{trace}((I - \pi) \text{Var}(y)) + \mathbb{E}(y^T) (I - \pi) \mathbb{E}(y) = \\ &= \text{trace}(I - \pi) \sigma^2 + (X^T (I - \pi) X) \beta = \sigma^2 \text{trace}(I - \pi) + \beta^T X^T (I - \pi) X \beta = \\ &= \sigma^2 (\text{trace}(I) - \text{trace}(\pi)) + \beta^T X^T (I - \pi) X \beta = (n - 1) \sigma^2 + \beta^T X^T (I - \pi) X \beta.\end{aligned}$$

Проделаем аналогичные вычисления для $\mathbb{E}(ESS)$

$$\begin{aligned}\mathbb{E}(ESS) &= \mathbb{E}(y^T (H - \pi) y) = \text{trace}((H - \pi) \text{Var}(y)) + \mathbb{E}(y^T) (H - \pi) \mathbb{E}(y) = \\ &= \text{trace}(H - \pi) \sigma^2 + (X^T (H - \pi) X) \beta = \sigma^2 \text{trace}(H - \pi) + \beta^T X^T (H - \pi) X \beta = \\ &= \sigma^2 (\text{trace}(H) - \text{trace}(\pi)) + \beta^T X^T (H - \pi) X \beta = (k - 1) \sigma^2 + \beta^T X^T (H - \pi) X \beta\end{aligned}$$

и для $\mathbb{E}(RSS)$

$$\mathbb{E}(RSS) = \mathbb{E}(y^T (I - H) y) = \text{trace}((I - H) \text{Var}(y)) + \mathbb{E}(y^T) (I - H) \mathbb{E}(y) =$$

$$\begin{aligned}
&= \text{trace}(I - H)\sigma^2 + (X\beta)^T(I - H)(X\beta) = \sigma^2 \text{trace}(I - H) + \beta^T X^T(I - H)X\beta = \\
&= \sigma^2(\text{trace}(I) - \text{trace}(H)) + \beta^T X^T(I - H)X\beta = (n - k)\sigma^2 + \beta^T X^T(I - H)X\beta = (n - k)\sigma^2.
\end{aligned}$$

Таким образом, получаем

$$\mathbb{E}(TSS) = (n - 1)\sigma^2 + \beta^T X^T(I - \pi)X\beta,$$

$$\mathbb{E}(ESS) = (k - 1)\sigma^2 + \beta^T X^T(H - \pi)X\beta,$$

$$\mathbb{E}(RSS) = (n - k)\sigma^2.$$

Задача 29. (Hansen 4.14)

Задана модель $y = X\beta + u$, для которой выполняются предпосылки теоремы Гаусса — Маркова. Вас интересует величина $\theta = \beta^2$. Получены МНК-оценки коэффициентов: $\hat{\beta}$, $V_{\hat{\beta}} = \mathbb{V}\text{ar}[\hat{\beta} \mid X]$. Кажется, неплохой идеей будет оценить θ как $\hat{\theta} = \hat{\beta}^2$.

- Найдите $\mathbb{E}[\hat{\theta} \mid X]$. Является ли $\hat{\theta}$ смещённой?
- Предложите способ коррекции смещения для получения несмещённой оценки $\hat{\theta}^*$, используя результаты предыдущего пункта.

Решение.

Здесь нужны решения

Задача 30. Рассмотрим модель регрессии $y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i$. Все предпосылки теоремы Гаусса — Маркова выполнены. Дополнительно предположим, что $u_i \sim \mathcal{N}(0, \sigma^2)$, $i = 1, \dots, n$. Дополнительно известно, что на самом деле $\beta_2 = \dots = \beta_k = 0$.

- Найдите $\mathbb{E}(R^2)$.

Решение:

Модель без ограничений:

$$y_i = \beta_1 + \beta_2 x_{i1} + \dots + \beta_k x_{ik} + u_i.$$

Модель с ограничениями (истинная модель!):

$$y_i = \beta_1 + u_i.$$

Тогда F-статистика имеет следующий вид:

$$F = \frac{R^2/(k - 1)}{(1 - R^2)/(n - k)} \sim F(k - 1, n - k).$$

Выразим R^2 :

$$R^2(n - k) = F(1 - R^2)(n - k)$$

Факт дня №1: Если $X \sim F(k_1, k_2)$, то $Y = \frac{\frac{k_1}{k_2} X}{1 + \frac{k_1}{k_2} X} \sim \text{Beta}\left(\frac{k_1}{2}, \frac{k_2}{2}\right)$.

Используя факт дня №1, получаем:

$$R^2 = \frac{(k-1)F}{(n-k) + (k-1)F} = \frac{\frac{k-1}{n-k}F}{1 + \frac{k-1}{n-k}F} \sim \text{Beta} \left(\frac{k-1}{2}, \frac{n-k}{2} \right).$$

Тогда чтобы посчитать математическое ожидание R^2 , надо вспомнить, чему равно математическое ожидание для $\text{Beta} \left(\frac{k-1}{2}, \frac{n-k}{2} \right)$:

$$E(R^2) = \frac{\frac{k-1}{2}}{\frac{k-1}{2} + \frac{n-k}{2}} = \frac{k-1}{n-1}.$$

Что нам даёт полученный результат? Математическое ожидание коэффициента детерминации линейно по k . То есть даже при включении в модель лишних факторов R^2 все равно продолжает линейно расти!

б) Найдите $\mathbb{E}(R_{\text{adj}}^2)$.

Решение:

Скорректированный коэффициент детерминации имеет вид:

$$R_{\text{adj}}^2 = 1 - (1 - R^2) \frac{n-1}{n-k}.$$

Рассчитаем математическое ожидание:

$$\begin{aligned} E(R_{\text{adj}}^2) &= E \left(1 - (1 - R^2) \frac{n-1}{n-k} \right) = 1 - \frac{n-1}{n-k} + \frac{n-1}{n-k} E(R^2) = \\ &= 1 - \frac{n-1}{n-k} + \frac{n-1}{n-k} \frac{k-1}{n-1} = 0. \end{aligned}$$

Скорректированный R^2 помог решить проблему линейного роста по k !

Задача 31. У овечки Долли был набор данных из n наблюдений для которого были выполнены предпосылки теоремы Гаусса — Маркова. Овечка Долли клонировала каждое наблюдение по одному разу и дописала каждое наблюдение-клон сразу после исходного наблюдения.

- а) Как выглядит ковариационная матрица ошибок для нового набора данных?
- б) Как изменится ответ на (а), если Долли клонирует только последнее наблюдение n раз?

Решение. а) Ковариационная матрица будет содержать блоки B на диагонали

$$\mathbb{V}\text{ar}(u) = \begin{pmatrix} B & 0 & 0 & \dots \\ 0 & B & 0 & \dots \\ 0 & 0 & B & \dots \\ \dots & & & \end{pmatrix},$$

где каждый блок равен $B = \begin{pmatrix} \sigma^2 & \sigma^2 \\ \sigma^2 & \sigma^2 \end{pmatrix}$.

- б) Ковариационная матрица будет состоять из четырех блоков: два блока нулевые, левый верхний блок пропорционален единичной матрицы, а все элементы правого нижнего блока равны σ^2 :

$$\mathbb{V}\text{ar}(u) = \begin{pmatrix} \sigma^2 \cdot I & 0 \\ 0 & S \end{pmatrix},$$

где I — единичная матрица, а все $S_{ij} = \sigma^2$.

6.9. Компьютерные задачи для семинара:

Генерация R2 для вывода распределения

Генерация смещения

Генерация лишних регрессоров

Реальный пример с лишним регрессорами (тип знаки зодиака и ретроградный)

Какая-то длинная задача, которую из темы в тему и в ней находить потом нарушения предпосылок?

<https://colab.research.google.com/drive/1wFrLyGcVVETx96jS93I4z8asgAQwqIdw?usp=sharing>

6.10. Домашнее задание:

6.11. Чёрный трэк:

Умножение блочных матриц. Если размеры блоков допускают операцию умножения, то:

$$\left[\begin{array}{c|c} A & B \\ \hline C & D \end{array} \right] \cdot \left[\begin{array}{c|c} E & F \\ \hline G & H \end{array} \right] = \left[\begin{array}{c|c} AE + BG & AF + BH \\ \hline CE + DG & CF + DH \end{array} \right].$$

Формула Фробениуса (блочное обращение).

$$\left[\begin{array}{c|c} A & B \\ \hline C & D \end{array} \right]^{-1} = \left[\begin{array}{c|c} A^{-1} + A^{-1}BH^{-1}CA^{-1} & -A^{-1}BH^{-1} \\ \hline -H^{-1}CA^{-1} & H^{-1} \end{array} \right],$$

где A — невырожденная квадратная матрица размерности $[n \times n]$, D — квадратная матрица размерности $[k \times k]$, $H = D - CA^{-1}B$.

Задача 32. Пусть истинной является модель $y = X_1\beta_1 + X_2\beta_2 + u$, где X_1, X_2 — матрицы признаков размерностей $[n \times k_1]$ и $[n \times k_2]$ соответственно. Вместо истинной модели вы оцениваете модель вида $y = X_1\beta_1 + v$, где v — вектор случайной ошибки, удовлетворяющий предпосылкам теоремы Гаусса — Маркова.

- а) Будет ли МНК-оценка вектора параметров β_1 несмещённой?
- б) Будет ли несмещённой МНК-оценка дисперсии случайной ошибки?
- в) Рассчитайте $\text{Var}(\hat{\beta}_1)$. Не противоречит ли полученный результат теореме Гаусса — Маркова?

Задача 33. Пусть истинной является модель $y = X_1\beta_1 + u$, где X_1 — матрица признаков размерности $[n \times k_1]$. Вместо истинной модели вы оцениваете модель вида $y = X_1\beta_1 + X_2\beta_2 + v$, где X_2 — матрица признаков размерности $[n \times k_2]$, v — вектор случайной ошибки, удовлетворяющий предпосылкам теоремы Гаусса — Маркова.

- а) Будет ли МНК-оценка вектора параметров β_1 несмещённой?
- б) Будет ли несмещённой МНК-оценка дисперсии случайной ошибки?
- в) Рассчитайте $\text{Var}(\hat{\beta}_1)$. Не противоречит ли полученный результат теореме Гаусса — Маркова?

7. Доверительные интервалы для коэффициентов

Построение доверительных интервалов для МНК оценок. Проверка гипотез. Асимптотика без нормальности ошибок. Нормальность ошибок.

7.1. Случай многомерного нормального распределения

сопроводить оценкой правдоподобия и показать, что она совпадает с МНК

Напомним несколько фактов про многомерное нормальное распределение.

Начнём с классического определения:

Определение 7.1 (многомерное нормальное распределение). Вектор v имеет многомерное невырожденное нормальное распределение, $v \sim \mathcal{N}(\mu, C)$, если его совместная функция плотности равна

$$f(v) = (2\pi)^{-n/2} \det(C)^{-1/2} \exp\left(-\frac{1}{2}(v - \mu)^T C^{-1}(v - \mu)\right),$$

где n — размерность вектора v .

Заметим, что совместный закон распределения нормального вектора v полностью определён его ожиданием $\mathbb{E}(v)$ и его ковариационной матрицей $\text{Var}(v)$. Никакие другие параметры в совместную функцию плотности не входят.

Для многомерного нормального распределения нет разницы между независимостью и некоррелированностью:

Теорема 7.2 (некоррелированность и независимость для нормального вектора). Если нормальный вектор v состоит из двух подвекторов, $v = (x, y)$, то $\text{Cov}(x, y) = 0$ если и только если подвекторы x и y независимы.

Доказательство. Докажем в одну сторону. Если подвекторы x и y независимы, то $\text{Cov}(x, y) = 0$. А теперь изящно докажем в обратную сторону. Если $\text{Cov}(x, y) = 0$, то вся ковариационная матрица $\text{Var}(v)$ ровно такая же как и в случае независимых x и y . Остаётся лишь вспомнить, что $\mathbb{E}(v)$ и $\text{Var}(v)$ полностью определяют закон распределения нормального вектора v , а значит компоненты обязаны быть независимы. \square

Также для многомерного нормального распределения нет разницы между условным ожиданием $\mathbb{E}(y \mid x)$ и наилучшим линейным приближением $\text{BestLin}(y \mid x)$, другими словами функция $\mathbb{E}(y \mid x)$ линейна по x .

Задача 34. Рассмотрим совместное нормальное распределение

$$\begin{pmatrix} x \\ y \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} C_{xx} & C_{xy} \\ C_{yx} & C_{yy} \end{pmatrix}\right)$$

- а) Найдите наилучшее линейное приближение $\text{BestLin}(y \mid x)$.
- б) Найдите условное ожидание $\mathbb{E}(y \mid x)$.
- в) Найдите условную дисперсию $\text{Var}(y \mid x)$.

Решение.

Здесь рассказать про определение bestlin в векторном случае?

- а) Пусть $\text{BestLin}(y \mid x) = a + Bx$. Мы хотим, чтобы ошибка линейной аппроксимации $r = y - \text{BestLin}(y \mid x)$ была некоррелирована с x ,

$$\mathbb{Cov}(y - \text{BestLin}(y \mid x), x) = 0.$$

Другими словами,

$$\mathbb{Cov}(y, x) = \mathbb{Cov}(\text{BestLin}(y \mid x), x) = 0.$$

Подставим $\text{BestLin}(y \mid x) = a + Bx$.

$$\mathbb{Cov}(y, x) = \mathbb{Cov}(a + Bx, x) = \mathbb{Cov}(Bx, x) = B \mathbb{Cov}(x, x) = B \mathbb{Var}(x).$$

Отсюда $C_{yx} = BC_{xx}$ и $B = C_{yx}C_{xx}^{-1}$. Кроме того, ошибка линейной аппроксимации должна иметь нулевое ожидание, следовательно,

$$\mathbb{E}(y) = \mathbb{E}(\text{BestLin}(y \mid x)) = a + B \mathbb{E}(x).$$

Получаем уравнение на a :

$$\mu_y = a + C_{yx}C_{xx}^{-1}\mu_x$$

И ответ,

$$\begin{cases} B = C_{yx}C_{xx}^{-1} \\ a = \mu_y - C_{yx}C_{xx}^{-1}\mu_x \end{cases}$$

- б) Для нормально распределённой пары векторов нулевая ковариация равносильная независимости. Следовательно, ошибка аппроксимации $r = y - \text{BestLin}(y \mid x)$ и x независимы. Отсюда мы получаем, что для многомерного нормально распределённого вектора (x, y)

$$\mathbb{E}(y \mid x) = \text{BestLin}(y \mid x) = a + Bx$$

- в) Условная дисперсия — это безусловная дисперсия ошибки прогноза,

$$\mathbb{Var}(y \mid x) = \mathbb{Var}(a + Bx + r \mid x) = \mathbb{Var}(r \mid x) = \mathbb{Var}(r).$$

Осталось вспомнить, что $y = a + Bx + r$, прогноз $a + Bx$ и ошибка r некоррелированы,

$$\mathbb{Var}(y) = \mathbb{Var}(a + Bx) + \mathbb{Var}(r).$$

Значит,

$$\mathbb{Var}(y \mid x) = \mathbb{Var}(r) = \mathbb{Var}(y) - B \mathbb{Var}(x) B^T = C_{yy} - C_{yx}C_{xx}^{-1}C_{xx}C_{xx}^{-1}C_{xy} = C_{yy} - C_{yx}C_{xx}^{-1}C_{xy}.$$

Отметим, что для компонент x и y нормального вектора (x, y) условная дисперсия получилась постоянной и не зависящей от x .

Для ненормального распределения условное ожидание $\mathbb{E}(y \mid x)$ и условная дисперсия $\mathbb{Var}(y \mid x)$ вполне могут быть нелинейными.

Введём дополнительную предпосылку $(u \mid X) \sim \mathcal{N}(0, \sigma^2 I)$. Учитывая, что $\hat{\beta} = (X^T X)^{-1} X^T y = \beta + (X^T X)^{-1} X^T u$, получаем

$$(\hat{\beta} \mid X) \sim \mathcal{N}(\beta, \sigma^2 (X^T X)^{-1}).$$

7.2. Независимость оценок β и $\hat{\sigma}^2$

МНК-оценка вектора коэффициентов β имеет вид

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

Несмещённая оценка дисперсии случайной ошибки:

$$\hat{\sigma}^2 = \frac{RSS}{n - k} = \frac{\hat{u}^T \hat{u}}{n - k}.$$

Распишем

$$\hat{\beta} = (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X + u) = \beta + (X^T X)^{-1} X^T u = \beta + Au,$$

где $A = (X^T X)^{-1} X^T$.

В случае, когда случайный вектор ошибок u является нормальным, можно показать, что оценки $\hat{\beta}$ и \hat{u} будут независимыми.

При $u \sim \mathcal{N}(0, \sigma^2 I_n)$ случайные векторы $\hat{\beta}$ и \hat{u} имеют совместное многомерное нормальное распределение. Покажем, что $\hat{\beta}$ и \hat{u} являются некоррелированными, из чего следует, что они также будут и независимыми, что справедливо для нормально распределённых векторов:

$$\text{Cov}(\hat{\beta}, \hat{u}) = \text{Cov}(\beta + Au, Mu) = \text{Cov}(Au, Mu) = AM \text{Var}(u) = \sigma^2 AM = 0, \text{ так как } AM = 0.$$

Так как $\hat{\sigma}^2$ есть функция от случайного вектора \hat{u} , то оценки $\hat{\beta}$ и $\hat{\sigma}^2$ также независимы.

7.3. Проверка гипотез о параметрах

$$H_0 : \beta_j = \beta_j^0$$

$$H_1 : \beta_j \neq \beta_j^0$$

$$t = \frac{\hat{\beta}_j - \beta_j^0}{\text{se}(\hat{\beta}_j)} \stackrel{H_0}{\sim} t(n - k)$$

Проверка гипотезы о незначимости модели в целом

$$H_0 : \beta_2 = \dots = \beta_k = 0$$

$$H_1 : \sum_{j=2}^k \beta_j^2 > 0$$

$$F = \frac{R^2 / (k - 1)}{(1 - R^2) / (n - k)} \stackrel{H_0}{\sim} F(k - 1, n - k).$$

8. Бутстрэп

Бутстрэп. Классический бутстрэп до регрессии и бутстрэп в регрессии. Метод наименьших модулей. Чёрный трек: возможно, разные варианты бутстрэпа в регрессии? ВСА-бутстрэп до регрессии?

9. Выбор функциональной формы

Дамми-переменные и их интерпретация. Функциональные формы: полиномы, логарифмы, интерпретация коэффициентов. Информационные критерии.

Чёрный трек: Структурные сдвиги. Тест Чоу. Локально-линейная регрессия (LOESS).

10. Гетероскедастичность

Гетероскедастичность. Тестирование гетероскедастичности. Робастные оценки. Доступный обобщённый МНК.

Задачи для доски:

Хансен: во сколько раз может быть недооценена дисперсия из-за гетероскедастичности

Коммент: акцент на робастных ошибках, тестирование и обобщённый МНК — кратко.

11. Мультиколлинеарность и метод главных компонент

Мультиколлинеарность и метод главных компонент.

Чёрный трек: несколько взглядов на метод главных компонент? LASSO?

12. Эндогенность

Эндогенность. Инструментальные переменные. Ошибка измерения регрессора. Двухшаговый МНК.

13. Эффекты воздействия

Оценка эффектов воздействия. ATE. LATE. Четкий (sharp) и нечеткий (fuzzy) разрывный регрессионный дизайн (RDD).

Чёрный трек: Метод разность разностей (DiD). Динамический метода разность разностей (Event Study).

14. Задачи

15. Логистическая регрессия: точечные оценки

Логистическая регрессия: Бинарный и упорядоченный логит. Точечные оценки, прогнозы. Интерпретация предельных эффектов.

Чёрный трек: Множественные логиты. Неупорядоченные, условные, смешанные логиты.

16. Логистическая регрессия: доверительные интервал

Логистическая регрессия: доверительные интервалы и проверка гипотез.

Чёрный трек: разные хоббиты

16.1. Смещение, цензурирование и ■■■■■■

Представим себе ситуацию, в которой зависимая количественная не всегда наблюдаема. Для моделирования этой ситуации мы введём скрытую латентная переменная y_i^* , которая линейно зависит от предиктора x_i , как обычно,

$$y_i^* = x_i^T \beta + u_i, \quad y^* = X^T \beta + u$$

Бинарная переменная $z_i \in \{0, 1\}$ равна 1 в случае, если мы наблюдаем y_i^* .

Возможно несколько случаев:

	наблюдаемость y^*	наблюдаемость x	наблюдаемость
Цензурирование censored model	зависит от y^*	всегда	
Усечение truncated model	зависит от y^*	если наблюдаем y^*	
Выборочное смещение sample selection	зависит от w	всегда	всегда
Переключающиеся режимы switching regimes	всегда, w переключает тип зависимости	всегда	всегда

Представим себе, что мы открыли дорогой ресторан. К нам заглядывают клиенты. Часть клиентов ужасаются от ценника и убегают, $y_i^* < 0$. Часть клиентов остаются и ужинают у нас, $y_i^* > 0$. Вместо нуля можно выбрать другой порог, но с нулём чуть-чуть удобнее.

16.2. Цензурирование

Рассмотрим самый распространённый вариант цензурирования: вместо отрицательных значений латентной переменной y_i^* мы видим нули.

Эта модель известна как тобит модель типа I, type I Tobit model.

$$\begin{cases} y_i^* = x_i^T \beta + u_i, & y^* = X^T \beta + u \\ (u \mid X) \sim \mathcal{N}(0, \sigma^2 I) \\ y_i = \max\{y_i^*, 0\} \\ (x_i, y_i) \text{ наблюдаемы при любых } i \end{cases}$$

Лог-функция правдоподобия равна

$$\ell(\beta, \sigma) = \sum_{y_i=0} \ln F(-x_i^T \beta / \sigma) + \sum_{y_i>0} \ln f((y_i - x_i^T \beta) / \sigma) - \sum_{y_i>0} \ln \sigma$$

16.3. Усечение

$$\begin{cases} y_i^* = x_i^T \beta + u_i, & y^* = X^T \beta + u \\ (u \mid X) \sim \mathcal{N}(0, \sigma^2 I) \\ y_i = \max\{y_i^*, 0\} \\ (x_i, y_i) \text{ наблюдаемы, если } y_i > 0 \end{cases}$$

Лог-функция правдоподобия равна

$$\ell(\beta, \sigma) = \sum_{y_i > 0} \ln f((y_i - x_i^T \beta)/\sigma) - \sum_{y_i > 0} \ln F(x_i^T \beta/\sigma) - \sum_{y_i > 0} \ln \sigma$$

16.4. Три осмысленных условных ожидания

Ожидание латентной переменной показывает, сколько в среднем планирует потратить гость ресторана на ужин, ещё не видевший цен, полезность от ужина,

$$m^*(x_i) = \mathbb{E}(y_i^* | x_i) = x_i^T \beta$$

Предельный эффект для латентной переменной

$$\partial \mathbb{E}(y_i^* | x_{ij}) / \partial x_{ij} = \beta_j$$

Ожидание цензурированной переменной, $y_i = \max\{y_i^*, 0\}$, сколько в среднем потратит человек, заглянувший в ресторан, с учётом того, что часть уйдёт испугавшись ценника

$$m(x_i) = \mathbb{E}(y_i | x_i) = x_i^T \beta F(x_i^T \beta/\sigma) + \sigma f(x_i^T \beta/\sigma)$$

Предельный эффект для цензурированной переменной

$$\partial \mathbb{E}(y_i | x_{ij}) / \partial x_{ij} =$$

Условное ожидание усечённой переменной, $(y_i | y_i^* > 0)$, средний чек в ресторане

$$m^\#(x_i) = \mathbb{E}(y_i | x_i, y_i^* > 0) = x_i^T \beta + \sigma \text{IMR}(x_i^T \beta/\sigma),$$

где $\text{IMR}(s)$ — обратное отношение Миллса, inverse Mills ratio,

$$\text{IMR}(s) = \mathbb{E}(v | v + s > 0) = f(s)/F(s), \quad v \sim \mathcal{N}(0; 1)$$

Предельный эффект для ожидания усечённой переменной

Выборочное смещение

Переключающиеся режимы

List of Theorems

2.1	Определение (наилучшее линейное приближение)	4
2.2	Определение (линейно-независимые случайные величины)	5
5.1	Определение (матрица Якоби)	12
5.2	Определение (матрица-шляпница)	17
5.3	Теорема	19
5.4	Определение (коэффициент детерминации)	20
5.5	Определение (коэффициент детерминации)	20
5.6	Определение (матрица-шляпница)	20
5.7	Теорема	21
5.8	Теорема (Теорема Фриша — Во — Ловелла (англ. Frisch–Waugh–Lovell theorem, FWL theorem))	22

5.9	Определение (LOOCV)	24
5.10	Теорема (связь обычных и кросс-валидационных остатков)	24
6.1	Определение (наилучшая линейная аппроксимация)	37
6.2	Определение (линейная независимость)	37
6.3	Определение (независимость случайных величин)	37
6.4	Теорема (Гаусс — Марков для парной линейной регрессии)	40
6.5	Теорема (Гаусс — Марков)	46
6.6	Определение (положительно полуопределённая форма)	48
6.7	Теорема (свойства положительно полуопределённой матрицы)	48
6.8	Теорема (ожидание суммы квадратов остатков)	49
7.1	Определение (многомерное нормальное распределение)	62
7.2	Теорема (некоррелированность и независимость для нормального вектора)	62
[PP12]	K. B. Petersen и M. S. Pedersen. The Matrix Cookbook. 2012. URL: http://www2.compute.dtu.dk/pubdb/pubs/3274-full.html . Матричные тождества, матричные распределения... Всё, что вы хотели узнать о матрицах, но боялись спросить! Обратите внимание на фамилии авторов :)	
[GD18]	Olya Gnailova и Boris Demeshev. How Gauss and Markov met Pythagoras: geometry in econometrics. 2018. URL: https://github.com/olyagnailova/gauss-markov-pythagoras . Как встретились Гаусс, Марков и Пифагор? Куча прикольных геометрических фактов и интерпретаций!	
