

Лекции по математической статистике

Шабанов Д.А., конспектировал Алексей Хачиянц

Содержание

1 Лекция 1	3
1.1 Ради чего мы собрались?	3
1.2 Точечное оценивание	3
1.3 Метод моментов	5
1.4 Выборочные квантили	6
2 Лекция 2	10
2.1 Сравнение оценок	10
2.2 Неравенство Рао-Крамера	13
2.3 Информация Фишера	15
2.4 Многомерное неравенство Рао-Крамера	17
3 Лекция 3	18
3.1 Оценки максимального правдоподобия	18
4 Лекция 4	24
4.1 Теорема Бахадура и следствия из неё	24
4.2 Байесовский подход	30
5 Лекция 5	31
5.1 Байесовские оценки	31
5.2 Минимаксные оценки	33
6 Лекция 6	36
7 Лекция 7	37
7.1 Достаточные статистики и оптимальные оценки	37
7.2 Критерий факторизации Неймана-Фишера	40
8 Лекция 8	42
8.1 Критерий факторизации Неймана-Фишера: продолжение	42
8.2 Полнота	44
9 Лекция 9	47
9.1 Линейная регрессия	47
9.2 Гауссовская линейная регрессия	50
10 Лекция 10	51
10.1 Доверительные интервалы и области	51
10.2 Проверка гипотез	54
11 Лекция 11	55
11.1 Сравнение критериев	55
12 Лекция 12	58
12.1 Критерий согласия	58
12.2 Параметрический хи-квадрат	60
13 Лекция 13	60
14 Лекция 14	62

15 Лекция 15	62
15.1 Слабая сходимость	62
15.2 Случайные процессы	64
15.3 Назад к критерию Колмогорова	65
16 Лекция 16	66
17 Лекция 17	66
17.1 Ранговые методы	66

1 Лекция 1

1.1 Ради чего мы собрались?

Допустим, что у нас есть какое-то наблюдение \mathbf{X} , то есть какой-то случайный вектор (хотя в асимптотических случаях будем считать, что \mathbf{X} бесконечномерен). Далее, распределение вектора \mathbf{X} считается неизвестным. Основная задача математической статистики состоит в том, чтобы *выдать обоснованное мнение* о распределении \mathbf{X} . Рассмотрим несколько основных подзадач:

- (а) Допустим, что мы знаем, что искомое распределение P принадлежит какому-то классу распределений \mathcal{P} . Мы хотим проверить, можно ли выделить более узкий подкласс распределений $\mathcal{P}_0 \subset \mathcal{P}$ такой, что $P \in \mathcal{P}_0$. Данный вид задач называется *проверкой гипотез*.

Пример: допустим, что нам известно, что случайный вектор X имеет нормальное распределение и мы хотим проверить гипотезу о том, что X имеет нормальное распределение с нулевым матожиданием.

- (б) Теперь предположим, что распределение P пришло из параметрического семейства: $P \in \{P_\theta : \theta \in \Theta\}$ и мы хотим оценить истинное значение параметра θ_0 , то есть построить оценку истинного значения θ_0 . Такие задачи называются *точечным оцениванием* (если мы оцениваем каким-то значением) или же *интервальным оцениванием* (если мы в качестве оценки предлагаем какую-то область).

- (с) Пусть наблюдение содержит пары: $\mathbf{X} = ((X_1, Y_1), \dots, (X_n, Y_n))$. Возникает вопрос: можно ли сказать, что Y_i и X_i независимы?

Пример: Допустим, что мы собираем данные о цвете волос и цвете глаз у людей и пытаемся проверить, есть ли зависимость между ними.

- (d) Допустим, что наблюдение разбилось на две части: $\mathbf{X} = (Y_1, \dots, Y_m) \sqcup (Z_1, \dots, Z_n)$. Вопрос таков: можно ли сказать, что Y_i и Z_j равны по распределению? Другими словами, можно ли сказать, что в наблюдении действительно будут одинаково распределённые случайные величины?

1.2 Точечное оценивание

Начнём с точечного оценивания. Пусть \mathbf{X} — выборка (случайный вектор с независимыми и одинаково распределёнными компонентами) из неизвестного параметрически заданного распределения $P \in \{P_\theta : \theta \in \Theta\}$, причём параметром является набор из k действительных чисел: $\Theta \subseteq \mathbb{R}^k$, $k \geq 1$. По выборке мы должны каким-то образом оценить истинное значение параметра θ . Для этого строятся *оценки*. Дадим определение:

Определение 1. *Оценка* — это борелевская¹ функция $T: \mathcal{X} \mapsto \Theta$, где \mathcal{X} — выборочное пространство. Проще говоря, это функция от выборки.

Для оценок нет единого обозначения. Если оценивается параметр θ , то оценку обозначают через $\hat{\theta}(\mathbf{X})$, $\theta^*(\mathbf{X})$ или же $\tilde{\theta}(\mathbf{X})$ (и так далее, но суть ясна).

Мы научились делать какие-то предсказания насчёт значения θ . Но как проверить адекватность предсказания? Для этого нужно выделить какие-то полезные свойства оценок и проверять их. Выпишем четыре основных свойства.

Определение 2. Пусть $\hat{\theta}(\mathbf{X})$ — это оценка параметра θ . Тогда будем называть оценку $\hat{\theta}$ *несмещённой*, если для любого $\theta \in \Theta$

$$E_\theta[\hat{\theta}(\mathbf{X})] = \theta.$$

В данном случае E_θ означает, что при взятии математического ожидания мы предполагаем, что выборка \mathbf{X} взята из распределения P_θ .

Определение 3. Пусть размер выборки \mathbf{X} увеличивается и $\hat{\theta}_n(\mathbf{X}) = \hat{\theta}_n(X_1, \dots, X_n)$. Тогда последовательность оценок $\hat{\theta}_n(\mathbf{X})$ будет называться *состоятельной* оценкой параметра θ , если для любого $\theta \in \Theta$

$$\hat{\theta}_n(\mathbf{X}) \xrightarrow{P_\theta} \theta.$$

Понятие *сильно состоятельной* оценки совпадает с определением состоятельной оценки с тем отличием, что в ней сходимость по вероятности заменяется на сходимость почти наверное:

$$\hat{\theta}_n(\mathbf{X}) \xrightarrow{P_\theta\text{-п.н.}} \theta.$$

¹Напомним, что борелевская функция — это отображение такое, что для него полный прообраз борелевского множества будет борелевским множеством.

Определение 4. Пусть размер выборки \mathbf{X} увеличивается и $\hat{\theta}_n(\mathbf{X}) = \hat{\theta}_n(X_1, \dots, X_n)$. Тогда последовательность оценок $\hat{\theta}_n(\mathbf{X})$ будет называться *асимптотически нормальной* оценкой параметра θ , если для любого $\theta \in \Theta$

$$\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \xrightarrow{d_\theta} \mathcal{N}(\mathbf{0}, \Sigma(\theta)),$$

где $\Sigma(\theta)$ — это *асимптотическая дисперсия* $\hat{\theta}_n(\mathbf{X})$.

Теперь приведём общую идею того, как можно строить оценки. Допустим, что мы смогли найти функционал G такой, что для всех $\theta \in \Theta$ $G(\mathbf{P}_\theta) = \theta$. Тогда «хорошей» оценкой можно считать $\hat{\theta}(\mathbf{X}) = G(\mathbf{P}_n^*)$, где \mathbf{P}_n^* — это *эмпирическое распределение*:²

$$\mathbf{P}_n^*(B) = \frac{1}{n} \sum_{k=1}^n [x_k \in B]$$

На эмпирическое распределение можно смотреть, как на дискретное распределение, равновероятно сосредоточенное в точках выборки \mathbf{X} . Теперь рассмотрим несколько примеров:

(а) Допустим, что функционал G устроен по следующему правилу:

$$G(\mathbf{P}_\theta) = \int_{\mathcal{X}} g(x) \mathbf{P}_\theta(dx).$$

Это интеграл Лебега по мере \mathbf{P}_θ от функции g . Людям, не знакомым с ним, можно читать его следующим образом: если \mathbf{P}_θ — это абсолютно непрерывная вероятностная мера с плотностью p_θ , то

$$\int_{\mathcal{X}} g(x) \mathbf{P}_\theta(dx) = \int_{\mathcal{X}} g(x) p_\theta(x) dx.$$

Если же \mathbf{P}_θ есть дискретная вероятностная мера, то

$$\int_{\mathcal{X}} g(x) \mathbf{P}_\theta(dx) = \sum_{x \in \mathcal{X}} g(x) \mathbf{P}_\theta(\{x\}).$$

В таком случае оценка будет иметь вид

$$G(\mathbf{P}_n^*) = \frac{1}{n} \sum_{i=1}^n g(X_i) \equiv \overline{g(\mathbf{X})}.$$

$\overline{g(\mathbf{X})}$ обычно называют *выборочным средним*.

(б) Теперь скажем, что оценка имеет вид линейной комбинации функций от порядковых статистик (такие оценки обычно называют *L-оценками*):

$$\hat{\theta}_n(\mathbf{X}) = \sum_{i=1}^n \alpha_i \varphi_i(X_{(i)}).$$

На всякий случай напомним определение порядковой статистики. Пусть $\omega \in \Omega$ и $x_k = X_k(\omega)$. Далее, перенумеруем последовательность так, чтобы она шла по возрастанию: $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$. Эту последовательность называют *вариационным рядом*. Случайную величину $X_{(k)}(\omega) = x_{(k)}$ называют *k-й порядковой статистикой*. Сразу же заметим, что $X_{(1)} = \min\{X_1, \dots, X_n\}$ и $X_{(n)} = \max\{X_1, \dots, X_n\}$.

В качестве примера возьмём G такой, что $G(\mathbf{P}_\theta)$ будет равен α -квантилю \mathbf{P}_θ . Напомню, что α -квантиль c_α функции распределения F равен минимальному значению x , в котором $F(x) \geq \alpha$:

$$c_\alpha = \min\{x : F(x) \geq \alpha\}.$$

Тогда *выборочным α -квантилем* назовём

$$G(\mathbf{P}_n^*) = \begin{cases} X_{(n\alpha)}, & n\alpha \text{ целое} \\ X_{(\lceil n\alpha \rceil)}, & n\alpha \text{ не целое} \end{cases}$$

²Здесь и в дальнейшем для обозначения индикатора будем использовать *нотацию Айверсона*: $[P] = 1$, если P правда и 0 иначе.

(с) Теперь скажем, что функционал G устроен следующим образом:

$$G(Q) = \arg \min_{\theta \in \Theta} \int_{\mathcal{X}} \psi(x, \theta) Q(dx).$$

Такие оценки называют *М-оценками*. Тогда

$$G(P_n^*) = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{k=1}^n \psi(X_k, \theta).$$

В качестве примера М-оценки можно взять оценку максимального правдоподобия: в ней $\psi(x, \theta) = -\log p_\theta(x)$.

Попробуем разобраться, в каких задачах можно построить «хорошие» оценки в том смысле, что они будут обладать свойствами, введённых выше.

Кто может определять распределение? Первая мысль, которая приходит в голову — моменты. Из этого мы сразу получаем так называемый *метод моментов*. Но этот метод плох: у него есть много недостатков, а плюс один — на него можно давать задачи на контрольной. Основная сложность состоит в том, что нужно уметь считать моменты, как функции от параметра, что может быть совсем нетривиально или же невозможно (особенно в случае, когда распределение далеко не табличное).

Для примера рассмотрим распределение Коши $C(\theta, 1)$. Обычно его задают плотностью $f(x; \theta, 1)$, которая равна

$$f(x; \theta, 1) = \frac{1}{\pi(1 + (x - \theta)^2)}.$$

Если взять случайную величину с таким распределением, то у неё не будет матожидания, дисперсии да и всех моментов выше первого. Конечно, можно посчитать матожидание логарифма (в таком случае интеграл сойдётся), но выразить её, как явную функцию от θ , не получится. Небольшое примечание: в данном распределении θ выступает в качестве медианы.

Второй подход предлагает использовать квантили. Очевидно, что этот подход универсален, так как по любой функции распределения можно построить квантили. Но у этого метода есть сложность — квантили считать ещё сложнее, чем моменты. Так что этот метод в основном подходит только для определения «сдвига»: если плотность симметрична относительно какой-то точки, то эта точка служит медианой, а медиану на практике считать не так уж и сложно.

Ещё одна идея завязывается на М-оценках. Описать её явно не так уж и просто, но общую философию можно сказать так: «мы живём в наиболее вероятном мире». Допустим, мы посчитали функционал и определили θ при фиксированном \mathbf{X} , при котором он максимален (или минимален). Тогда эту θ берут в качестве настоящей, так как мы считаем, что она отражает действительность.

1.3 Метод моментов

Пусть \mathbf{X} — выборка из распределения $P \in \{P_\theta: \theta \in \Theta\}$, причём $\Theta \subseteq \mathbb{R}^k$. Далее, возьмём *пробные функции* $g_1(x), \dots, g_k(x)$ такие, что вектор

$$m(\theta) = (E_\theta[g_1(X_1)], \dots, E_\theta[g_k(X_1)])$$

задаёт биекцию между Θ и $m(\Theta)$. Тогда оценкой $\hat{\theta}(\mathbf{X})$ по методу моментов с пробными функциями $g_1(x), \dots, g_k(x)$ называется решение системы уравнений (относительно θ)

$$\begin{cases} m_1(\theta) = \overline{g_1(\mathbf{X})} \\ \dots \\ m_k(\theta) = \overline{g_k(\mathbf{X})} \end{cases}$$

Здесь $m_j(\theta) = E_\theta[g_j(X_1)]$. Так как $m(\theta)$ задаёт биекцию, то

$$\hat{\theta}(\mathbf{X}) = m^{-1}(\overline{g_1(\mathbf{X})}, \dots, \overline{g_k(\mathbf{X})}).$$

Обычно в качестве пробных функций берут x, x^2, \dots, x^k .

Оказывается, что с теоретической точки зрения оценки с помощью метода моментов не так уж и плохи.

Лемма. Пусть $\hat{\theta}_n(\mathbf{X})$ — это оценка θ по методу моментов. Тогда

- (a) Если m^{-1} непрерывна, то $\hat{\theta}_n(\mathbf{X})$ есть сильно состоятельная оценка для θ .
- (b) Если m^{-1} дифференцируема и $E_\theta[g_j^2(X_1)] < +\infty$ для всех $j = 1, \dots, k$, то $\hat{\theta}_n(\mathbf{X})$ есть асимптотически нормальная оценка θ .

Доказательство. Для начала заметим, что по усиленному закону больших чисел

$$\overline{g_j(\mathbf{X})} = \frac{1}{n} \sum_{k=1}^n g_j(X_k) \xrightarrow{P_{\theta-\text{п.н.}}} E_\theta[g_j(X_1)].$$

Следовательно, по теореме о наследовании сходимости

$$\hat{\theta}_n(\mathbf{X}) = m^{-1}(\overline{g_1(\mathbf{X})}, \dots, \overline{g_k(\mathbf{X})}) \xrightarrow{P_{\theta-\text{п.н.}}} m^{-1}(m_1(\theta), \dots, m_k(\theta)) = m^{-1}(m(\theta)) = \theta.$$

Для доказательства асимптотической нормальности заметим следующее. Введём следующие векторы: $\mathbf{Y}_i = (g_1(X_i), \dots, g_k(X_i))$. Далее, так как элементы выборки независимы и одинаково распределены, то \mathbf{Y}_i тоже независимы и одинаково распределены. Так как мы предполагаем, что $E_\theta[g_j^2(X_1)] < +\infty$, то матрица ковариаций Σ невырождена. Тогда мы можем применить многомерную центральную предельную теорему:

$$\sqrt{n} \left(\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} - E_\theta[\mathbf{Y}_1] \right) \xrightarrow{d_\theta} \mathcal{N}(\mathbf{0}, \Sigma(\theta))$$

Для дальнейших рассуждений нужна следующая

Лемма (о наследовании асимптотической сходимости). Пусть $\xi_i \in \mathbb{R}^m$ — последовательность случайных векторов, для которых существует вектор $\mathbf{a} \in \mathbb{R}^m$ и матрица $\Sigma \in \mathbb{R}^{m \times m}$ такие, что

$$\sqrt{n}(\xi_n - \mathbf{a}) \xrightarrow{d_\theta} \mathcal{N}(\mathbf{0}, \Sigma).$$

Далее, пусть $H: \mathbb{R}^m \mapsto \mathbb{R}^n$ — это дифференцируемая функция. Тогда

$$\sqrt{n}(H(\xi_n) - H(\mathbf{a})) \xrightarrow{d_\theta} H'(\mathbf{a})\mathcal{N}(\mathbf{0}, \Sigma),$$

где $H'(\mathbf{a})$ — это матрица Якоби функции H в точке \mathbf{a} .

Теперь применим эту лемму, положив $H = m^{-1}$:

$$\sqrt{n} \left(m^{-1} \left(\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} \right) - m^{-1}(E_\theta[\mathbf{Y}_1]) \right) = \sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \xrightarrow{d_\theta} (m^{-1})'(m(\theta))\mathcal{N}(\mathbf{0}, \Sigma(\theta)).$$

Осталось заметить, что $(m^{-1})'(m(\theta))\mathcal{N}(\mathbf{0}, \Sigma(\theta)) = \mathcal{N}(\mathbf{0}, \mathbf{A}(\theta))$ для некоторой матрицы $\mathbf{A}(\theta)$. □

1.4 Выборочные квантили

Теперь приступим ко второму методу описания функции распределения — через квантили. Оказывается, что у выборочных квантилей есть очень хорошее свойство.

Теорема 1. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из распределения $F(x)$ с плотностью $f(x)$, а z_p — p -квантиль распределения $F(x)$. Далее, пусть $f(x)$ непрерывно дифференцируема в окрестности z_p и $f(z_p) > 0$. Тогда выборочный квантиль есть асимптотически нормальная оценка теоретического квантиля:

$$\sqrt{n}(X_{[np]+1} - z_p) \xrightarrow{d_\theta} \mathcal{N} \left(0, \frac{p(1-p)}{f^2(z_p)} \right).$$

Доказательство. Начнём с того, что посчитаем плотность k -й порядковой статистики. Для этого заметим, что если $X_{(k)} \leq x$, то хотя бы k элементов выборки не больше x . Тогда, если $F_k(x)$ — это функция распределения k -й статистики, то

$$F_k(x) = P(X_{(k)} \leq x) = \sum_{s=k}^n \binom{n}{s} F^s(x) (1 - F(x))^{n-s}.$$

Теперь посчитаем плотность $f_k(x)$ функции распределения $F_k(x)$, продифференцировав её:

$$f_k(x) = \sum_{s=k}^n s f(x) \binom{n}{s} F^{s-1}(x) (1 - F(x))^{n-s} - \sum_{s=k}^{n-1} (n-s) f(x) \binom{n}{s} F^s(x) (1 - F(x))^{n-s-1}.$$

Осталось упростить эту сумму. Для этого заметим, что

$$s \binom{n}{s} = \frac{n!}{(s-1)!(n-s)!} = n \binom{n-1}{s-1}, \quad (n-s) \binom{n}{s} = \frac{n!}{s!(n-s-1)!} = n \binom{n-1}{s}.$$

Тогда, заменив индекс суммирования в первой сумме на $t = s - 1$, получим следующее:

$$f_k(x) = \sum_{t=k-1}^{n-1} n f(x) \binom{n-1}{t} F^t(x) (1-F(x))^{n-t-1} - \sum_{t=k}^{n-1} n f(x) \binom{n-1}{t} F^t(x) (1-F(x))^{n-t-1}.$$

Но тогда

$$f_k(x) = n f(x) \binom{n-1}{k-1} F^{k-1}(x) (1-F(x))^{n-k}.$$

Теперь покажем, что если мы возьмём последовательность случайных величин T_n , построенных по правилу

$$T_n = \frac{f(z_p) \sqrt{n}}{\sqrt{p(1-p)}} (X_{(k)} - z_p), \quad \text{где } k = [np] + 1,$$

то $T_n \xrightarrow{d_\theta} \mathcal{N}(0, 1)$. Для этого покажем один промежуточный факт: если ξ — случайная величина из распределения с плотностью f_ξ , а $\eta = a\xi + b$, где $a > 0$ и b — константы, то η имеет плотность f_η , равную

$$f_\eta(x) = \frac{1}{a} f_\xi \left(\frac{x-b}{a} \right).$$

Это несложно понять, если заметить следующее:

$$F_\eta(x) = P(\eta \leq x) = P(a\xi + b \leq x) = P\left(\xi \leq \frac{x-b}{a}\right) = F_\xi\left(\frac{x-b}{a}\right).$$

Дифференцируя по x , мы получим желаемое. Воспользуемся этим. Пусть $q_n(x)$ — плотность T_n . Тогда

$$q_n(x) = \frac{\sqrt{p(1-p)}}{f(z_p) \sqrt{n}} f_k \left(z_p + \frac{x \sqrt{p(1-p)}}{f(z_p) \sqrt{n}} \right).$$

Для удобства введём следующее обозначение:

$$t_n = z_p + \frac{x \sqrt{p(1-p)}}{f(z_p) \sqrt{n}}.$$

Тогда

$$q_n(x) = n \binom{n-1}{k-1} F(t_n)^{k-1} f(t_n) (1-F(t_n))^{n-k} \sqrt{\frac{p(1-p)}{n}} \frac{1}{f(z_p)}.$$

Теперь сделаем финт ушами и скажем, что

$$q_n(x) = A_1(n) A_2(n) A_3(n),$$

где

$$\begin{aligned} A_1(n) &= \frac{f(t_n)}{f(z_p)}, \\ A_2(n) &= n \binom{n-1}{k-1} \sqrt{\frac{p(1-p)}{n}} p^{k-1} (1-p)^{n-k}, \\ A_3(n) &= \left(\frac{F(t_n)}{p} \right)^{k-1} \left(\frac{1-F(t_n)}{1-p} \right)^{n-k}. \end{aligned}$$

Найдём пределы всех трёх выражений при $n \rightarrow \infty$:

1. В первом всё очень просто: достаточно заметить, что $t_n \rightarrow z_p$ при $n \rightarrow \infty$ и $f(x)$ непрерывна в окрестности z_p . Тогда

$$\frac{f(t_n)}{f(z_p)} \xrightarrow{n \rightarrow \infty} 1 \implies \lim_{n \rightarrow \infty} A_1(n) = 1.$$

2. Для этого пункта заметим, что

$$A_2(n) = k \binom{n}{k} \sqrt{\frac{p(1-p)}{n}} p^{k-1} (1-p)^{n-k}.$$

Далее, воспользуемся формулой Стирлинга:

$$A_2(n) \sim k \frac{\sqrt{2\pi n} \left(\frac{n}{e}\right)^n}{\sqrt{2\pi k} \left(\frac{k}{e}\right)^k \sqrt{2\pi(n-k)} \left(\frac{n-k}{e}\right)^k} \sqrt{\frac{p(1-p)}{n}} p^{k-1} (1-p)^{n-k}.$$

Теперь вспомним, что $k = [np] + 1$. Но из этого следует, что $k \sim np$ и $n - k \sim n(1-p)$. Тогда

$$A_2(n) \sim np \frac{\sqrt{2\pi n} \left(\frac{n}{e}\right)^n}{\sqrt{2\pi np} \left(\frac{k}{e}\right)^k \sqrt{2\pi n(1-p)} \left(\frac{n-k}{e}\right)^k} \sqrt{\frac{p(1-p)}{n}} p^{k-1} (1-p)^{n-k}.$$

Упростим:

$$A_2(n) \sim \frac{1}{\sqrt{2\pi}} \left(\frac{np}{k}\right)^k \left(\frac{n(1-p)}{n-k}\right)^{n-k}.$$

Теперь докажем, что $A_2(n) \rightarrow 1/\sqrt{2\pi}$ при $n \rightarrow \infty$. Заметим, что

$$A_2(n) \sim \frac{1}{\sqrt{2\pi}} \exp \left\{ k \ln \frac{np}{k} + (n-k) \ln \frac{n(1-p)}{n-k} \right\}.$$

Далее, k отличается от np не более, чем на 1. Тогда логарифмы можно разложить в ряд Тейлора в нуле:

$$\begin{aligned} \ln \frac{np}{k} &= \frac{np}{k} - 1 + O\left(\frac{1}{k^2}\right) = \frac{np-k}{k} + O\left(\frac{1}{n^2}\right), \\ \ln \frac{n(1-p)}{n-k} &= \frac{n(1-p)}{n-k} - 1 + O\left(\frac{1}{(n-k)^2}\right) = \frac{k-np}{n-k} + O\left(\frac{1}{n^2}\right). \end{aligned}$$

Следовательно,

$$A_2(n) \sim \frac{1}{\sqrt{2\pi}} \exp \left\{ np - k + k - np + O\left(\frac{1}{n^2}\right) \right\} \xrightarrow{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}}.$$

3. Для этого пункта нужно заметить, что $F(t_n) \rightarrow F(z_p) = p$ при $n \rightarrow \infty$. Из этого можно сделать вывод, что

$$\left(\frac{F(t_n)}{p}\right)^{k-1} \sim \left(\frac{F(t_n)}{p}\right)^k.$$

Далее, заметим, что

$$A_3(n) \sim \exp \left\{ k \ln \frac{F(t_n)}{p} + (n-k) \ln \frac{1-F(t_n)}{1-p} \right\}.$$

Теперь разложим $F(t_n)$ в ряд Тейлора в точке z_p , пользуясь непрерывной дифференцируемостью $f(x)$:

$$\begin{aligned} F(t_n) &= F(z_p) + (t_n - z_p)f(z_p) + \frac{1}{2}(t_n - z_p)^2 f'(z_p) + o((t_n - z_p)^2) \\ &= p + \sqrt{\frac{p(1-p)}{n}}x + \frac{f'(z_p)}{f^2(z_p)} \frac{p(1-p)}{2n}x^2 + o((t_n - z_p)^2). \end{aligned}$$

Тогда

$$\frac{F(t_n)}{p} = 1 + \sqrt{\frac{1-p}{np}}x + \frac{f'(z_p)}{f^2(z_p)} \frac{1-p}{2n}x^2 + o\left(\frac{1}{n}\right)$$

Следовательно, если взять от этого логарифм, то его можно разложить в ряд Тейлора:

$$\ln \frac{F(t_n)}{p} = \sqrt{\frac{1-p}{np}}x + \frac{f'(z_p)}{f^2(z_p)} \frac{1-p}{2n}x^2 + o\left(\frac{1}{n}\right) - \frac{1}{2} \frac{1-p}{np}x^2 + o\left(\frac{1}{n}\right).$$

Тогда

$$k \ln \frac{F(t_n)}{p} = \sqrt{\frac{1-p}{np}}xk + \frac{f'(z_p)}{f^2(z_p)} \frac{k(1-p)}{2n}x^2 - \frac{1}{2} \frac{1-p}{p} \frac{k}{n}x^2 + o(1).$$

Аналогичными рассуждениями можно показать, что

$$(n-k) \ln \frac{1-F(t_n)}{1-p} = -\sqrt{\frac{p}{n(1-p)}} x(n-k) - \frac{f'(z_p)}{f^2(z_p)} \frac{p(n-k)}{2n} x^2 - \frac{1}{2} \frac{p}{1-p} \frac{n-k}{n} x^2 + o(1).$$

Теперь воспользуемся тем, что $|k-np| \leq 1$. Это означает, что $k = np + O(1/n)$. Тогда

$$\begin{aligned} k \ln \frac{F(t_n)}{p} &= \sqrt{\frac{1-p}{np}} x \left(np + O\left(\frac{1}{n}\right) \right) + \frac{1}{2} \frac{f'(z_p)}{f^2(z_p)} \frac{1-p}{n} x^2 \left(np + O\left(\frac{1}{n}\right) \right) \\ &\quad - \frac{1}{2} \frac{1-p}{p} \frac{x^2}{n} \left(np + O\left(\frac{1}{n}\right) \right) + o(1) \\ &= x \sqrt{np(1-p)} + O\left(\frac{1}{\sqrt{n}}\right) + \frac{p(1-p)}{2} \frac{f'(z_p)}{f^2(z_p)} x^2 + O\left(\frac{1}{n^2}\right) \\ &\quad - \frac{1-p}{2} x^2 + O\left(\frac{1}{n^2}\right) + o(1) \\ &= x \sqrt{np(1-p)} + \frac{p(1-p)}{2} \frac{f'(z_p)}{f^2(z_p)} x^2 - \frac{1-p}{2} x^2 + o(1). \end{aligned}$$

Аналогично поступая со вторым слагаемым, получим, что

$$(n-k) \ln \frac{1-F(t_n)}{1-p} = -x \sqrt{np(1-p)} - \frac{p(1-p)}{2} \frac{f'(z_p)}{f^2(z_p)} x^2 - \frac{p}{2} x^2 + o(1).$$

Следовательно,

$$A_3(n) \sim \exp \left\{ -\frac{x^2}{2} + o(1) \right\} \xrightarrow{n \rightarrow \infty} \exp \left\{ -\frac{x^2}{2} \right\}.$$

В итоге мы получили, что для всех x

$$\lim_{n \rightarrow \infty} q_n(x) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{x^2}{2} \right\}.$$

Это означает, что $q_n(x)$ будет равномерно сходиться к плотности $\mathcal{N}(0, 1)$ на любом компакте. Из равномерной сходимости на отрезке $[a, b]$ следует, что

$$\lim_{n \rightarrow \infty} (F_{T_n}(b) - F_{T_n}(a)) = \lim_{n \rightarrow \infty} \int_a^b q_n(x) dx = \Phi(b) - \Phi(a).$$

Теперь нужно доказать, что $F_{T_n}(x) \rightarrow \Phi(x)$ для всех x , где $\Phi(x)$ — это стандартное нормальное распределение. Заметим, что для любого $a < b$

$$\overline{\lim}_{n \rightarrow \infty} |F_{T_n}(b) - \Phi(b)| \leq \overline{\lim}_{n \rightarrow \infty} |F_{T_n}(b) - F_{T_n}(a) + \Phi(b) - \Phi(a)| + \overline{\lim}_{n \rightarrow \infty} |F_{T_n}(a) - \Phi(a)|$$

Далее, первый предел уходит в ноль. Тогда

$$\overline{\lim}_{n \rightarrow \infty} |F_{T_n}(b) - \Phi(b)| \leq \overline{\lim}_{n \rightarrow \infty} |F_{T_n}(a) - \Phi(a)|.$$

Однако

$$\overline{\lim}_{n \rightarrow \infty} |F_{T_n}(a) - \Phi(a)| \leq \overline{\lim}_{n \rightarrow \infty} (F_{T_n}(a) + \Phi(a)) = \Phi(a) + \overline{\lim}_{n \rightarrow \infty} F_{T_n}(a).$$

Далее, снова сделаем финт ушами, пользуясь тем, что $F_{T_n}(a) < 1$:

$$\Phi(a) + \overline{\lim}_{n \rightarrow \infty} F_{T_n}(a) \leq \Phi(a) + \overline{\lim}_{n \rightarrow \infty} (F_{T_n}(a) - F_{T_n}(-a) + 1) = \Phi(a) + \Phi(a) - \Phi(-a) + 1.$$

Но эту сумму можно сделать сколь угодно малой, устремив a к $-\infty$. Следовательно, $T_n \xrightarrow{d_\theta} \mathcal{N}(0, 1)$ и мы получаем желаемое. \square

Посмотрим на частный случай квантили, а именно на медиану. Стоит заметить, что выборочную медиану определяют не как выборочный $1/2$ -квантиль:

Определение 5. Выборочной медианой $\hat{\mu}$ выборки X_1, \dots, X_n называют величину

$$\hat{\mu}(\mathbf{X}) = \begin{cases} X_{(k+1)}, & n = 2k + 1 \\ (X_{(k)} + X_{(k+1)})/2, & n = 2k \end{cases}$$

Оказывается, что для такого определения свойство асимптотической нормальности тоже выполняется:
Упражнение. В условиях теоремы про асимптотическую нормальность выборочной квантили

$$\sqrt{n}(\hat{\mu}(\mathbf{X}) - z_{1/2}) \xrightarrow{d_\theta} \mathcal{N}\left(0, \frac{1}{4f^2(z_{1/2})}\right).$$

В качестве примера рассмотрим выборку \mathbf{X} из распределения Коши $\mathcal{C}(\theta, 1)$. У данного распределения всё очень плохо с моментами, но при этом асимптотическая нормальность выборочной медианы выполняется:

$$\sqrt{n}(\hat{\mu}(\mathbf{X}) - \theta) \xrightarrow{d_\theta} \mathcal{N}\left(0, \frac{\pi^2}{4}\right).$$

2 Лекция 2

2.1 Сравнение оценок

Начнём с такой темы, как сравнение оценок. Очевидно, что оценок можно придумать навалом, можно и много оценок с хорошими свойствами наподобие состоятельности или асимптотической нормальности и прочими. Вопрос: кто из них лучше и как сравнивать их?

Пусть $\hat{\theta}(\mathbf{X})$ и $\theta^*(\mathbf{X})$ — оценки неизвестного параметра θ с «хорошими» свойствами. Грубо говоря, если оценка не состоятельна, то она нам не интересна глобально. Возникает вопрос: как сравнивать оценки $\hat{\theta}(\mathbf{X})$ и $\theta^*(\mathbf{X})$?

Базовый принцип сравнения завязан на так называемых функциях потерь и риска.

Определение 6. Функцию $\rho(x, y) \geq 0$ называют *функцией потерь*.

Простыми словами, функция потерь считает, сколько мы потеряем, если вместо x подставим y . Данная функция обычно должна иметь некоторые хорошие свойства (если она метрика — то это вообще прекрасно, но это не обязательно).

Приведём несколько примеров:

- (a) Самый базовый пример — это квадратичная функция потерь: $\rho(x, y) = (x - y)^2$ ($x, y \in \mathbb{R}$).
- (b) $\rho(x, y) = |x - y|$ — L_1 -метрика.
- (c) $\rho(\mathbf{x}, \mathbf{y}) = \langle \mathbf{A}(\mathbf{x} - \mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$, где $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{A} \succ 0$ (это означает, что матрица положительно определена).

Определение 7. *Функцией риска* оценки $\hat{\theta}(\mathbf{X})$ параметра θ называется средний размер потерь нашей оценки:

$$R(\theta, \hat{\theta}(\mathbf{X})) = \mathbb{E}_\theta[\rho(\theta, \hat{\theta}(\mathbf{X}))].$$

Как мы можем сравнивать оценки с помощью функции риска? Для этого есть аж четыре подхода.

1. Самый банальный подход называется *равномерным подходом*. Суть его крайне проста: просто сравниваем функции риска.

Определение 8. Пусть $\hat{\theta}(\mathbf{X})$ и $\theta^*(\mathbf{X})$ — оценки θ . Будем говорить, что оценка $\hat{\theta}(\mathbf{X})$ лучше, чем оценка $\theta^*(\mathbf{X})$ (в равномерном смысле), если для всех $\theta \in \Theta$ функция риска для $\hat{\theta}(\mathbf{X})$ не превосходит функцию риска для $\theta^*(\mathbf{X})$:

$$R(\theta, \hat{\theta}(\mathbf{X})) \leq R(\theta, \theta^*(\mathbf{X})),$$

причём существует параметр $\theta_0 \in \Theta$ такой, что в нём неравенство будет строгим: $R(\theta_0, \hat{\theta}(\mathbf{X})) < R(\theta_0, \theta^*(\mathbf{X}))$.

Как правило, если мы взяли хорошую функцию потерь (наподобие квадратичной), то равенство функций риска везде означает равенство оценок, так что последнее требование не совсем обязательно.

Что теперь напрашивается? Поиск лучшей оценки. Будем называть оценку $\hat{\theta}(\mathbf{X})$ наилучшей, если она лучше любой другой в равномерном подходе. Однако поиск наилучшей оценки — это достаточно бессмысленная задача, так как мы её не найдём. Рассмотрим следующий пример. Пусть $\rho(x, y) = (x - y)^2$ — квадратичная функция потерь, параметр $\theta \in [0, 1]$, а оценка возвращает тождественный ноль: $\hat{\theta}(\mathbf{X}) \equiv 0$. Тогда функция риска будет равна

$$R(\theta, \hat{\theta}(\mathbf{X})) = \mathbb{E}_\theta[(\theta - \hat{\theta}(\mathbf{X}))^2] = \theta^2.$$

Однако, если $\theta = 0$, то $R(\theta, \hat{\theta}(\mathbf{X})) = 0$. Аналогично, если оценка возвращает не 0, а какое-то $\theta_0 \in [0, 1]$, то $R(\theta_0, \hat{\theta}(\mathbf{X})) = 0$. Но тогда если бы у нас была наилучшая оценка $\theta^*(\mathbf{X})$, то у неё должен быть тождественно нулевой риск: $R(\theta, \theta^*(\mathbf{X})) \equiv 0$. Тем самым мы получаем, что $\hat{\theta}(\mathbf{X}) = \theta$. Следовательно, наилучшей оценки в равномерном смысле просто нет.

Как тогда быть? Можно сказать, что будем смотреть не на все оценки, а только лишь на какой-то класс. Тогда можно ввести следующее определение:

Определение 9. Оценка $\hat{\theta}(\mathbf{X})$ называется наилучшей оценкой в классе оценок \mathcal{K} , если она лучше, чем любая другая оценка $\theta^*(\mathbf{X}) \in \mathcal{K}$. Если при подсчёте функции риска используется квадратичная функция потерь, то оценку $\hat{\theta}(\mathbf{X})$ называют *оптимальной* в классе оценок \mathcal{K} .

В качестве примера класса оценок можно сказать, что \mathcal{K} — это все несмещённые оценки $\tau(\theta)$.

2. *Байесовский подход.* Смысл этого подхода состоит в следующем: равномерный подход пытается считать, что все параметры в каком-то смысле равноправны и мы должны иметь неравенство для каждого $\theta \in \Theta$. Байесовский подход утверждает, что параметры не равноправны и некоторые параметры более вероятны, чем другие и мы пытаемся рассматривать параметр θ как случайную величину. Тогда нам важно выполнение неравенства на главных значениях, а на маловероятных уже без разницы. Формализуем это:

Определение 10. Пусть Q — некоторое вероятностное распределение на Θ с плотностью $q(t)$, а $\hat{\theta}(\mathbf{X})$ и $\theta^*(\mathbf{X})$ есть оценки θ . Тогда будем говорить, что оценка $\hat{\theta}(\mathbf{X})$ лучше оценки $\theta^*(\mathbf{X})$, если

$$R_Q(\hat{\theta}(\mathbf{X})) < R_Q(\theta^*(\mathbf{X})), \text{ где } R_Q(\hat{\theta}(\mathbf{X})) = \int_{\Theta} R(\theta, \hat{\theta}(\mathbf{X})) q(\theta) d\theta.$$

Далее, будем называть оценку $\hat{\theta}(\mathbf{X})$ *наилучшей*, если для неё R_Q минимально:

$$R_Q(\hat{\theta}(\mathbf{X})) = \inf_{\theta^*(\mathbf{X})} R_Q(\theta^*(\mathbf{X})).$$

Стоит заметить, что наилучшая оценка зависит от выбора распределения Q , то есть байесовский подход зависит от выбора распределения. По сути, на Q можно смотреть, как на априорное представление о том, какие значения могут принимать параметры. Грубо говоря, мы можем пытаться делать какие-то такие вещи: если какая-то оценка оказалась ближе к истине, чем другие, но их мы откидывать не хотим, тогда просто возьмём их с какими-то весами.

В качестве классического примера можно взять задачу про однорукого бандита.³ Допустим, что есть игровой автомат, у которого несколько ручек — для красоты будем считать, что ручек две, и вероятность выигрыша у одной ручки больше, чем у другой. Наша задача — оптимизировать стратегию среднего выигрыша. Самый простой подход состоит в следующем: прокрутим обе ручки по 100 раз и посмотрим, с какой ручки больше выигрыша (пусть вышло так, что с первой). Равномерный подход говорит, что первая ручка лучше и достаточно крутить только её. Байесовский же подход утверждает, что могло случиться так, что вторая ручка лучше, но мы её просто не использовали слишком малое количество раз и нам просто не повезло — если бы взяли миллион раз, то было бы лучше. Поэтому стратегия байесовского подхода будет в таком стиле: десять раз крутим первую ручку и один раз вторую.

3. *Минимаксный подход.* Этот подход достаточно тесно связан с байесовским подходом. Он апеллирует к следующему: мы вообще не хотим больших значений функций потерь, то есть если она мала для всех значений $\theta \in \Theta$, то нас это устраивает.

Определение 11. Пусть $\hat{\theta}(\mathbf{X})$ и $\theta^*(\mathbf{X})$ — оценки θ . Будем говорить, что оценка $\hat{\theta}(\mathbf{X})$ лучше, чем оценка $\theta^*(\mathbf{X})$ (в минимаксном смысле), если

$$\sup_{\theta \in \Theta} R(\theta, \hat{\theta}(\mathbf{X})) < \sup_{\theta \in \Theta} R(\theta, \theta^*(\mathbf{X})).$$

Оценку $\hat{\theta}(\mathbf{X})$ будем называть наилучшей, если для неё выполнено следующее равенство:

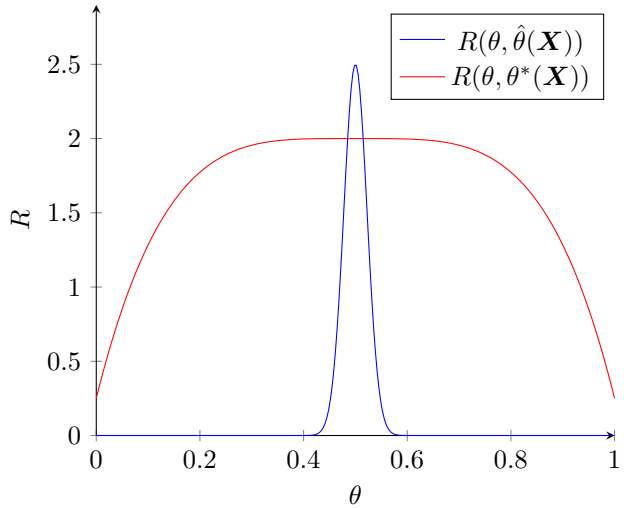
$$\sup_{\theta \in \Theta} R(\theta, \hat{\theta}(\mathbf{X})) = \inf_{\theta^*(\mathbf{X})} \sup_{\theta \in \Theta} R(\theta, \theta^*(\mathbf{X})).$$

³На самом деле это задача так называемого *последовательного анализа*.

Пусть у нас есть две оценки $\hat{\theta}(\mathbf{X})$ и $\theta^*(\mathbf{X})$ для параметра $\theta \in [0, 1]$ со следующими функциями риска (см. рисунок справа). Что можно сказать про эти оценки?

- Если в байесовском подходе сказать, что мы берём равномерное распределение, то есть $Q = U(0, 1)$, то оценка $\hat{\theta}(\mathbf{X})$ будет лучше оценки $\theta^*(\mathbf{X})$, так как у неё меньше площадь под кривой.
- Теперь рассмотрим минимаксный подход. В нём оценка $\theta^*(\mathbf{X})$ будет лучше оценки $\hat{\theta}(\mathbf{X})$, так как у неё меньше максимальное значение.
- В равномерном же подходе эти оценки несравнимы, так как есть и участки, где $R(\theta, \hat{\theta}(\mathbf{X})) > R(\theta, \theta^*(\mathbf{X}))$, так и участки, где $R(\theta, \hat{\theta}(\mathbf{X})) < R(\theta, \theta^*(\mathbf{X}))$.

Функции риска для $\hat{\theta}(\mathbf{X})$ и $\theta^*(\mathbf{X})$



Стоит заметить, что если в байесовском подходе взять распределение, сосредоточенное в центре, то оценка $\theta^*(\mathbf{X})$ будет лучше, чем $\hat{\theta}(\mathbf{X})$.

Всё это сравнение оценок приводит к понятию *допустимой оценки*. Суть в том, что равномерный подход наиболее сильный — если оценка лучше в равномерном подходе, то она будет лучше и в байесовском, и в минимаксном подходе. Поэтому имеет смысл смотреть только те оценки, для которых нет более хороших оценок в равномерном смысле, так как иначе можно просто взять более хорошую.

Определение 12. Будем называть оценку $\hat{\theta}(\mathbf{X})$ *допустимой*, если не существует оценки $\theta^*(\mathbf{X})$ такой, что $\theta^*(\mathbf{X})$ лучше, чем $\hat{\theta}(\mathbf{X})$ в равномерном подходе.

Есть ещё один подход для сравнения оценок, но он принципиально отличается от предыдущих тем, что он не использует функции риска и он работает только для асимптотически нормальных оценок.

4. *Асимптотический подход.* Он завязан на определении асимптотической дисперсии.

Определение 13. Пусть $\hat{\theta}_n(\mathbf{X})$ и $\theta_n^*(\mathbf{X})$ — асимптотически нормальные оценки параметра θ с асимптотическими дисперсиями $\sigma_1^2(\theta)$ и $\sigma_2^2(\theta)$ соответственно. В таком случае будем говорить, что оценка $\hat{\theta}_n(\mathbf{X})$ лучше, чем оценка $\theta_n^*(\mathbf{X})$, если для всех $\theta \in \Theta$ $\sigma_1^2(\theta) \leq \sigma_2^2(\theta)$.

Интуитивно на это можно смотреть, как на равномерное ограничение предельного среднеквадратичного уклонения. Почему это может быть важно? С помощью асимптотической дисперсии можно строить доверительные интервалы, и при её уменьшении будет уменьшаться длина интервала, что хорошо.

Для примера попробуем сравнить какие-нибудь оценки.

Задача 1. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из $\mathcal{N}(\theta, 1)$, $\bar{\mathbf{X}}$ — выборочное среднее, а $\hat{\mu}(\mathbf{X})$ — выборочная медиана. Сравните эти оценки в асимптотическом подходе.

Решение. Для начала сразу же уточним, что для нормального распределения математическое ожидание является медианой, поэтому оценивается одна и та же величина и задача действительно осмысленна. Далее, согласно центральной предельной теореме выборочное среднее является асимптотически нормальной оценкой:

$$\sqrt{n}(\bar{\mathbf{X}} - \theta) \xrightarrow{d_\theta} \mathcal{N}(0, D_\theta[X_1]) = \mathcal{N}(0, 1).$$

Заметим, что плотность $f(x)$ нормального распределения непрерывно дифференцируема и положительна. Следовательно, можно воспользоваться теоремой об асимптотической нормальности выборочной медианы:

$$\sqrt{n}(\hat{\mu}(\mathbf{X}) - \theta) \xrightarrow[n \rightarrow \infty]{d_\theta} \mathcal{N}\left(0, \frac{1}{4f^2(\theta)}\right) = \mathcal{N}\left(0, \frac{\pi}{2}\right).$$

Но тогда выборочная медиана хуже в асимптотическом подходе, чем выборочное среднее. И это неспроста. \square

Этот факт кажется не самым очевидным, так как известно, что медиана более устойчива к выбросам в выборке.

2.2 Неравенство Рао-Крамера

Будем считать, что мы работаем в классе \mathcal{K} несмещённых оценок $\tau(\theta)$. Поставим следующую задачу: как найти оптимальную оценку в этом классе? Другими словами, нам нужно решить следующую задачу: если $\hat{\theta}(\mathbf{X})$ — это несмещённая оценка $\tau(\theta)$, то нам нужно равномерно по всем $\theta \in \Theta$ минимизировать $E_{\theta}[(\hat{\theta}(\mathbf{X}) - \tau(\theta))^2] = D_{\theta}[\hat{\theta}(\mathbf{X})]$.

Оказывается, что если наложить так называемые *условия регулярности*, то можно предоставить нижнюю оценку для дисперсии оценки. Сформулируем эти условия:

1. Пусть множество параметров Θ — это открытый интервал на \mathbb{R} .
2. Далее, пусть параметрическое семейство распределений $\{P_{\theta} : \theta \in \Theta\}$ является доминируемым с плотностью $p_{\theta}(x)$. Немного разъясним это требование.

Определение 14. Параметрическое семейство распределений называется *доминируемым*, если все распределения в нём принадлежат к одному типу (абсолютно непрерывному или дискретному).

Если все распределения абсолютно непрерывны, то под плотностью $p_{\theta}(x)$ подразумевается обычная плотность вероятностного распределения P_{θ} . Если же все распределения дискретны, то $p_{\theta}(x) = P_{\theta}(\{x\})$.

В таком случае удобно ввести понятие *правдоподобия*. Это просто плотность выборки \mathbf{X} и она равна

$$p_{\theta}(\mathbf{X}) = \prod_{i=1}^n p_{\theta}(X_i).$$

3. Пусть $A = \{x : p_{\theta}(x) > 0\}$. Тогда A не зависит от θ .

Сделаем небольшое отступление: для того, чтобы доказывать некоторые утверждения одновременно и для дискретных, и для абсолютно непрерывных распределений, скажем, что⁴

$$\int_A p_{\theta}(x) \mu(dx) = \begin{cases} \int_A p_{\theta}(x) dx, & P_{\theta} \text{ абсолютно непрерывна} \\ \sum_{x \in A} p_{\theta}(x), & P_{\theta} \text{ дискретна} \end{cases}$$

4. Для любой статистики $S(\mathbf{X})$ такой, что $E_{\theta}[S^2(\mathbf{X})] < +\infty$, есть возможность дифференцирования под знаком интеграла:

$$\frac{\partial}{\partial \theta} E_{\theta}[S(\mathbf{X})] = E_{\theta}[S(\mathbf{X}) U_{\theta}(\mathbf{X})], \text{ где } U_{\theta}(\mathbf{X}) = \frac{\partial}{\partial \theta} \ln p_{\theta}(\mathbf{X}).$$

Теперь покажем, почему это то же самое, что и дифференцирование под знаком интеграла:

$$\begin{aligned} \frac{\partial}{\partial \theta} E_{\theta}[S(\mathbf{X})] &= \frac{\partial}{\partial \theta} \int_A S(x) p_{\theta}(x) \mu(dx) = \int_A S(x) \frac{\partial p_{\theta}(x)}{\partial \theta} \mu(dx) \\ &= \int_A S(x) \left(\frac{1}{p_{\theta}(x)} \frac{\partial p_{\theta}(x)}{\partial \theta} \right) p_{\theta}(x) \mu(dx) \\ &= \int_A S(x) \frac{\partial \ln p_{\theta}(x)}{\partial \theta} p_{\theta}(x) \mu(dx) = E_{\theta}[S(\mathbf{X}) U_{\theta}(\mathbf{X})]. \end{aligned}$$

Функцию $U_{\theta}(\mathbf{X})$ называют *вкладом* наблюдения \mathbf{X} .

5. Величина

$$I_{\mathbf{X}}(\theta) = E_{\theta} \left[\left(\frac{\partial}{\partial \theta} \ln p_{\theta}(\mathbf{X}) \right)^2 \right],$$

называемая *информацией* (по Фишеру) выборки \mathbf{X} , должна быть положительна и конечна для любых $\theta \in \Theta$.

Теперь можно сформулировать теорему, которая будет давать оценку снизу.

Теорема 2 (Неравенство Рао-Крамера). Пусть τ — это дифференцируемая функция от параметра. Далее, пусть $\hat{\theta}(\mathbf{X})$ — это несмещённая оценка $\tau(\theta)$ такая, что $E_{\theta}[\hat{\theta}^2(\mathbf{X})] < +\infty$ и выполнены условия регулярности. Тогда для всех $\theta \in \Theta$

$$D_{\theta}[\hat{\theta}(\mathbf{X})] \geq \frac{(\tau'(\theta))^2}{I_{\mathbf{X}}(\theta)}.$$

⁴В принципе, это связано с ремаркой из первой лекции, ибо здесь под μ подразумевается мера Лебега.

Доказательство. Воспользуемся четвёртым условием регулярности, подставив в неё $S(\mathbf{X}) \equiv 1$:

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathbb{E}_\theta[S(\mathbf{X})] &= \frac{\partial}{\partial \theta} \mathbb{E}_\theta[1] = 0, \\ \mathbb{E}_\theta[S(\mathbf{X})U_\theta(\mathbf{X})] &= \mathbb{E}_\theta[U_\theta(\mathbf{X})].\end{aligned}$$

Следовательно, $\mathbb{E}_\theta[U_\theta(\mathbf{X})] = 0$.

Далее, снова воспользуемся четвёртым свойством регулярности, подставив в него $S(\mathbf{X}) = \hat{\theta}(\mathbf{X})$:

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathbb{E}_\theta[S(\mathbf{X})] &= \frac{\partial}{\partial \theta} \mathbb{E}_\theta[\hat{\theta}(\mathbf{X})] = \tau'(\theta), \\ \mathbb{E}_\theta[S(\mathbf{X})U_\theta(\mathbf{X})] &= \mathbb{E}_\theta[\hat{\theta}(\mathbf{X})U_\theta(\mathbf{X})] = \mathbb{E}_\theta[(\hat{\theta}(\mathbf{X}) - \tau(\theta))U_\theta(\mathbf{X})]\end{aligned}$$

Тем самым $\tau'(\theta) = \mathbb{E}_\theta[(\hat{\theta}(\mathbf{X}) - \tau(\theta))U_\theta(\mathbf{X})]$. Теперь вспомним неравенство Коши-Буняковского-Шварца: для любых случайных величин X, Y с конечными вторыми моментами

$$\mathbb{E}^2[XY] \leq \mathbb{E}[X^2] \mathbb{E}[Y^2].$$

Тогда согласно этому неравенству

$$(\tau'(\theta))^2 \leq \mathbb{E}_\theta[(\hat{\theta}(\mathbf{X}) - \tau(\theta))^2] \mathbb{E}_\theta[U_\theta^2(\mathbf{X})] = D_\theta[\hat{\theta}(\mathbf{X})] I_\mathbf{X}(\theta).$$

Отсюда получаем желаемое:

$$D_\theta[\hat{\theta}(\mathbf{X})] \geq \frac{(\tau'(\theta))^2}{I_\mathbf{X}(\theta)}.$$

□

А можно ли вообще достичь эту оценку? В некоторых случаях можно.

Определение 15. Пусть $\hat{\theta}(\mathbf{X})$ — несмещённая оценка для $\tau(\theta)$. Будем называть её *эффективной*, если для неё достигается равенство в неравенстве Рао-Крамера, то есть для всех $\theta \in \Theta$

$$D_\theta[\hat{\theta}(\mathbf{X})] = \frac{(\tau'(\theta))^2}{I_\mathbf{X}(\theta)}.$$

Оказывается, для проверки оценки на эффективность даже есть критерий.

Теорема 3 (Критерий эффективности). *В условиях неравенства Рао-Крамера $\hat{\theta}(\mathbf{X})$ будет эффективной оценкой $\tau(\theta)$ тогда и только тогда, когда*

$$\hat{\theta}(\mathbf{X}) - \tau(\theta) = c(\theta)U_\theta(\mathbf{X}), \text{ где } c(\theta) = \frac{\tau'(\theta)}{I_\mathbf{X}(\theta)}.$$

Доказательство. Как известно, в неравенстве Коши-Буняковского-Шварца равенство достигается тогда и только тогда, когда между случайными величинами есть линейная зависимость. Следовательно, $\hat{\theta}(\mathbf{X})$ будет эффективной оценкой для $\tau(\theta)$ только в том случае, если существуют функции $c(\theta)$ и $a(\theta)$ такие, что

$$\hat{\theta}(\mathbf{X}) - \tau(\theta) = c(\theta)U_\theta(\mathbf{X}) + a(\theta).$$

Теперь посчитаем эти функции.

- Для начала возьмём матожидание:

$$\mathbb{E}_\theta[\hat{\theta}(\mathbf{X}) - \tau(\theta)] = \mathbb{E}_\theta[c(\theta)U_\theta(\mathbf{X}) + a(\theta)].$$

Тогда

$$\tau(\theta) - \tau(\theta) = c(\theta) \cdot 0 + a(\theta) \implies a(\theta) = 0.$$

- Теперь помножим обе части равенства на $U_\theta(\mathbf{X})$ и возьмём матожидание:

$$\mathbb{E}_\theta[(\hat{\theta}(\mathbf{X}) - \tau(\theta))U_\theta(\mathbf{X})] = \mathbb{E}_\theta[c(\theta)U_\theta^2(\mathbf{X})].$$

Но тогда

$$\tau'(\theta) = c(\theta)I_\mathbf{X}(\theta) \implies c(\theta) = \frac{\tau'(\theta)}{I_\mathbf{X}(\theta)}.$$

Тем самым получаем желаемое.

□

Теперь попробуем применить этот критерий.

Задача 2. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из распределения Бернулли $\text{Bin}(1, \theta)$, где $\theta \in (0, 1)$. Найдите эффективную оценку θ и информацию $I_{\mathbf{X}}(\theta)$.

Решение. Воспользуемся критерием эффективности. Для этого нужно посчитать вклад. Начнём с того, что запишем плотность выборки.

Небольшое отступление: если нам известно, что случайная величина ξ принимает значения x_1, \dots, x_m с вероятностями p_1, \dots, p_m , то плотность можно записать следующим образом:

$$p_{\xi}(x) = \prod_{i=1}^m p_i^{[x=x_i]}.$$

Воспользуемся этим и запишем правдоподобие:

$$p_{\theta}(\mathbf{X}) = \prod_{i=1}^n p_{\theta}(X_i) = \prod_{i=1}^n \theta^{[X_i=1]} (1-\theta)^{[X_i=0]} = \prod_{i=1}^n \theta^{X_i} (1-\theta)^{1-X_i}$$

Теперь несложно посчитать правдоподобие:

$$\begin{aligned} \ln p_{\theta}(\mathbf{X}) &= \sum_{k=1}^n X_k \ln \theta + (1 - X_k) \ln(1 - \theta) \\ U_{\theta}(\mathbf{X}) &= \frac{X_1 + \dots + X_n}{\theta} - \frac{n - X_1 - \dots - X_n}{1 - \theta} \\ &= n \left(\frac{\bar{X}}{\theta} - \frac{1 - \bar{X}}{1 - \theta} \right) = \frac{n}{\theta(1 - \theta)} (\bar{X} - \theta). \end{aligned}$$

Отсюда видно, что \bar{X} является эффективной оценкой θ и информация Фишера равна

$$I_{\mathbf{X}}(\theta) = \frac{n}{\theta(1 - \theta)}.$$

□

2.3 Информация Фишера

Начнём с того, что обобщим понятие информации по Фишеру:

Определение 16. Пусть $S(\mathbf{X})$ — это некоторая статистика с плотностью $g_{\theta}(x)$. Тогда информацией Фишера статистики $S(\mathbf{X})$ называется

$$I_S(\theta) = \mathbb{E}_{\theta} \left[\left(\frac{\partial}{\partial \theta} \ln g_{\theta}(S(\mathbf{X})) \right)^2 \right].$$

У информации Фишера есть несколько свойств, которые оправдывают её название.

1. Допустим, что распределение $S(\mathbf{X})$ не зависит от θ . Тогда $I_S(\theta) = 0$.

Доказательство. Это достаточно очевидное утверждение, так как в таком случае $\ln g_{\theta}(S(\mathbf{X}))$ не будет зависеть от θ и при дифференцировании по нему обратится в 0, а матожидание нуля есть ноль. □

2. Пусть $S(\mathbf{X})$ и $T(\mathbf{X})$ — независимые статистики для всех $\theta \in \Theta$ и выполнены условия регулярности. Тогда

$$I_{(S,T)}(\theta) = I_S(\theta) + I_T(\theta).$$

Доказательство. Пусть статистики $S(\mathbf{X})$ и $T(\mathbf{X})$ имеют плотности $f_{\theta}(s)$ и $g_{\theta}(t)$ соответственно. Тогда случайный вектор $(S(\mathbf{X}), T(\mathbf{X}))$ имеет совместную плотность $h_{\theta}(s, t) = f_{\theta}(s)g_{\theta}(t)$. Следовательно,

$$\ln h_{\theta}(s, t) = \ln f_{\theta}(s) + \ln g_{\theta}(t).$$

Из этого можно сделать вывод, что

$$\mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta} \ln h_{\theta}(S(\mathbf{X}), T(\mathbf{X})) \right] = \mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta} \ln f_{\theta}(S(\mathbf{X})) \right] + \mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta} \ln g_{\theta}(T(\mathbf{X})) \right] = 0.$$

Но тогда

$$\begin{aligned} I_{(S,T)}(\theta) &= \mathbb{D}_{\theta} \left[\frac{\partial}{\partial \theta} \ln h_{\theta}(S(\mathbf{X}), T(\mathbf{X})) \right] \\ &= \mathbb{D}_{\theta} \left[\frac{\partial}{\partial \theta} \ln f_{\theta}(S(\mathbf{X})) \right] + \mathbb{D}_{\theta} \left[\frac{\partial}{\partial \theta} \ln g_{\theta}(T(\mathbf{X})) \right] = I_S(\theta) + I_T(\theta). \end{aligned}$$

□

У этого свойства есть достаточно приятное следствие. Пусть

$$i(\theta) = \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \ln p_\theta(X_1) \right)^2 \right].$$

Тогда

$$I_{\mathbf{X}}(\theta) = \sum_{i=1}^n I_{X_i}(\theta) = n I_{X_1}(\theta) = n i(\theta).$$

3. Для любой статистики $S(\mathbf{X})$ $I_S(\theta) \leq I_{\mathbf{X}}(\theta)$.

Доказательство. Для начала покажем, что если у статистики $S(\mathbf{X})$ есть плотность $g_\theta(s)$, то

$$\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) \mid S(\mathbf{X}) \right] = \frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})).$$

Будем делать это по определению условного математического ожидания. Начнём с того, что функция справа является борелевской функцией от $S(\mathbf{X})$. Осталось проверить интегральное свойство: для любого борелевского множества B

$$\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) [S(\mathbf{X}) \in B] \right] = \mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) [S(\mathbf{X}) \in B] \right].$$

Заметим, что

$$\begin{aligned} \mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) [S(\mathbf{X}) \in B] \right] &= \int_A \frac{\partial \ln p_\theta(\mathbf{x})}{\partial \theta} p_\theta(\mathbf{x}) [S(\mathbf{x}) \in B] \mu(d\mathbf{x}) = \int_A \frac{\partial p_\theta(\mathbf{x})}{\partial \theta} [S(\mathbf{x}) \in B] \mu(d\mathbf{x}) \\ &= \frac{\partial}{\partial \theta} \int_A p_\theta(\mathbf{x}) [S(\mathbf{x}) \in B] \mu(d\mathbf{x}) = \frac{\partial}{\partial \theta} \mathbb{P}_\theta(S(\mathbf{X}) \in B) \end{aligned}$$

Но, с другой стороны

$$\begin{aligned} \frac{\partial}{\partial \theta} \mathbb{P}_\theta(S(\mathbf{X}) \in B) &= \frac{\partial}{\partial \theta} \int_B g_\theta(s) \mu'(ds) = \int_B \frac{\partial \ln p_\theta(s)}{\partial \theta} g_\theta(s) \mu'(ds) \\ &= \mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) [S(\mathbf{X}) \in B] \right]. \end{aligned}$$

Теперь будем смотреть на следующую величину:

$$M = \mathbb{E}_\theta \left[\frac{\partial \ln p_\theta(\mathbf{X})}{\partial \theta} \frac{\partial \ln g_\theta(S(\mathbf{X}))}{\partial \theta} \right].$$

Согласно неравенству Коши-Буняковского-Шварца

$$M^2 \leq \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) \right)^2 \right] \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) \right)^2 \right] = I_{\mathbf{X}}(\theta) I_S(\theta).$$

Теперь докажем, что $M = I_S(\theta)$. Для этого заметим, что

$$\mathbb{E}_\theta \left[\frac{\partial \ln p_\theta(\mathbf{X})}{\partial \theta} \frac{\partial \ln g_\theta(S(\mathbf{X}))}{\partial \theta} \right] = \mathbb{E}_\theta \left[\mathbb{E}_\theta \left[\frac{\partial \ln p_\theta(\mathbf{X})}{\partial \theta} \frac{\partial \ln g_\theta(S(\mathbf{X}))}{\partial \theta} \mid S(\mathbf{X}) \right] \right].$$

Далее, функцию от условия можно вынести из условного математического ожидания. Тогда

$$\begin{aligned} \mathbb{E}_\theta \left[\frac{\partial \ln p_\theta(\mathbf{X})}{\partial \theta} \frac{\partial \ln g_\theta(S(\mathbf{X}))}{\partial \theta} \right] &= \mathbb{E}_\theta \left[\frac{\partial \ln g_\theta(S(\mathbf{X}))}{\partial \theta} \mathbb{E}_\theta \left[\frac{\partial \ln p_\theta(\mathbf{X})}{\partial \theta} \mid S(\mathbf{X}) \right] \right] \\ &= \mathbb{E}_\theta \left[\left(\frac{\partial \ln g_\theta(S(\mathbf{X}))}{\partial \theta} \right)^2 \right] = I_S(\theta). \end{aligned}$$

Тогда

$$I_S^2(\theta) \leq I_{\mathbf{X}}(\theta) I_S(\theta) \implies I_S(\theta) \leq I_{\mathbf{X}}(\theta). \quad \square$$

4. Оценка $S(\mathbf{X})$ будет достаточной тогда и только тогда, когда для всех $\theta \in \Theta$ $I_S(\theta) = I_{\mathbf{X}}(\theta)$.

Доказательство. Это свойство будет доказано позднее. \square

Рассмотрим пример.

Задача 3. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из экспоненциального распределения $\text{Exp}(\theta)$. Далее, пусть $X_{(1)} = \min_{1 \leq i \leq n} X_i$. Найдите $I_{X_{(1)}}(\theta)$ и $I_{\mathbf{X}}(\theta)$.

Решение. Начнём с того, что найдём распределение $X_{(1)}$. Для этого заметим, что

$$P_\theta(X_{(1)} \leq x) = 1 - P_\theta(X_{(1)} \geq x) = 1 - \prod_{i=1}^n P_\theta(X_i \geq x) = 1 - e^{-n\theta x}.$$

Следовательно, $X_{(1)} \sim \text{Exp}(n\theta)$.

Теперь можно считать информации Фишера. Начнём с $I_{\mathbf{X}}(\theta)$. Вспомним, что $I_{\mathbf{X}}(\theta) = ni(\theta)$. Тогда

$$\frac{\partial}{\partial \theta} \ln p_\theta(X_1) = \frac{\partial}{\partial \theta} (\ln \theta - \theta X_1) = \frac{1}{\theta} - X_1.$$

Тогда

$$i(\theta) = E_\theta \left[\left(X_1 - \frac{1}{\theta} \right)^2 \right] = D_\theta[X_1] = \frac{1}{\theta^2} \implies I_{\mathbf{X}}(\theta) = \frac{n}{\theta^2}.$$

Теперь посчитаем $I_{X_{(1)}}(\theta)$ аналогичным образом. Заметим, что

$$\frac{\partial}{\partial \theta} \ln p_\theta(X_{(1)}) = \frac{\partial}{\partial \theta} (\ln n + \ln \theta - n\theta X_1) = \frac{1}{\theta} - nX_1.$$

Тогда

$$I_{X_{(1)}}(\theta) = E_\theta \left[\left(nX_1 - \frac{1}{\theta} \right)^2 \right] = n^2 D_\theta[X_{(1)}] = \frac{1}{\theta^2}. \quad \square$$

2.4 Многомерное неравенство Рао-Крамера

У неравенства Рао-Крамера есть многомерный аналог. Но для его доказательства нужно одно неравенство.

Теорема 4 (Матричное неравенство Коши-Буняковского-Шварца). Пусть Ψ и \mathbf{H} — случайные матрицы одного и того же размера и матрица $E_\theta[\mathbf{H}\mathbf{H}^\top]$ обратима. Тогда⁵

$$E_\theta[\Psi\Psi^\top] \succcurlyeq E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1} E_\theta[\mathbf{H}\Psi^\top],$$

причём равенство достигается тогда и только тогда, когда $\Psi = \mathbf{Z}\mathbf{H}$, где

$$\mathbf{Z} = E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1}.$$

Доказательство. Для начала заметим, что для любой матрицы \mathbf{A} матрица $\mathbf{A}\mathbf{A}^\top \succcurlyeq 0$. Тогда для любой матрицы \mathbf{Z} выполнено следующее:

$$(\Psi - \mathbf{Z}\mathbf{H})(\Psi - \mathbf{Z}\mathbf{H})^\top \succcurlyeq 0$$

Возьмём матожидание:

$$E_\theta[(\Psi - \mathbf{Z}\mathbf{H})(\Psi - \mathbf{Z}\mathbf{H})^\top] \succcurlyeq 0$$

Теперь раскроем матожидание по линейности:

$$E_\theta[\Psi\Psi^\top] - \mathbf{Z} E_\theta[\mathbf{H}\Psi^\top] - E_\theta[\Psi\mathbf{H}^\top]\mathbf{Z}^\top + \mathbf{Z} E_\theta[\mathbf{H}\mathbf{H}^\top]\mathbf{Z}^\top \succcurlyeq 0$$

Возьмём $\mathbf{Z} = E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1}$ и попробуем упростить выражение:

$$\begin{aligned} \mathbf{Z} E_\theta[\mathbf{H}\Psi^\top] &= E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1} E_\theta[\mathbf{H}\Psi^\top], \\ E_\theta[\Psi\mathbf{H}^\top]\mathbf{Z}^\top &= E_\theta[\Psi\mathbf{H}^\top](E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1})^\top \\ &= E_\theta[\Psi\mathbf{H}^\top]((E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1})^\top (E_\theta[\Psi\mathbf{H}^\top])^\top \\ &= E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1} E_\theta[\mathbf{H}\Psi^\top], \\ \mathbf{Z} E_\theta[\mathbf{H}\mathbf{H}^\top]\mathbf{Z}^\top &= E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1} E_\theta[\mathbf{H}\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1} E_\theta[\mathbf{H}\Psi^\top] \\ &= E_\theta[\Psi\mathbf{H}^\top](E_\theta[\mathbf{H}\mathbf{H}^\top])^{-1} E_\theta[\mathbf{H}\Psi^\top]. \end{aligned}$$

⁵Запись $\mathbf{A} \succcurlyeq \mathbf{B}$ означает, что матрица $\mathbf{A} - \mathbf{B}$ неотрицательно определена.

Следовательно,

$$\mathbb{E}_\theta[\Psi\Psi^\top] - \mathbb{E}_\theta[\Psi\mathbf{H}^\top](\mathbb{E}_\theta[\mathbf{H}\mathbf{H}^\top])^{-1}\mathbb{E}_\theta[\mathbf{H}\Psi^\top] \succcurlyeq 0.$$

Для критерия равенства нужно заметить, что

$$\mathbb{E}_\theta[(\Psi - \mathbf{Z}\mathbf{H})(\Psi - \mathbf{Z}\mathbf{H})^\top] = \mathbf{0} \iff \Psi = \mathbf{Z}\mathbf{H}.$$

□

Дело осталось за малым.

Теорема 5 (Многомерное неравенство Рао-Крамера). Пусть пространство параметров $\Theta \subseteq \mathbb{R}^k$, $k > 1$. Далее, пусть $\tau: \Theta \mapsto \mathbb{R}^k$ — это некоторая дифференцируемая по θ функция и её якобиан равен

$$\tau'(\theta) = \left\| \frac{\partial(\tau(\theta))_i}{\partial\theta_j} \right\|_{i,j=1}^k \in \mathbb{R}^{k \times k}.$$

Далее, $\hat{\theta}(\mathbf{X}) \in \mathbb{R}^k$ — это несмещённая оценка для $\tau(\theta)$ с матрицей ковариаций

$$\mathbf{D}_\theta[\hat{\theta}(\mathbf{X})] = \mathbb{E}_\theta[(\hat{\theta}(\mathbf{X}) - \tau(\theta))(\hat{\theta}(\mathbf{X}) - \tau(\theta))^\top]$$

Пусть $I_{\mathbf{X}}(\theta)$ — информация Фишера:

$$I_{\mathbf{X}}(\theta) = \left\| \mathbb{E}_\theta \left[\frac{\partial \ln p_\theta(\mathbf{X})}{\partial\theta_i} \frac{\partial \ln p_\theta(\mathbf{X})}{\partial\theta_j} \right] \right\|_{i,j=1}^k \in \mathbb{R}^{k \times k}$$

Другими словами, это матрица ковариаций вклада

$$U_\theta(\mathbf{X}) = \left(\frac{\partial}{\partial\theta_1} \ln p_\theta(\mathbf{X}), \dots, \frac{\partial}{\partial\theta_k} \ln p_\theta(\mathbf{X}) \right)^\top$$

Многомерное неравенство Рао-Крамера утверждает, что в условиях регулярности

$$\mathbf{D}_\theta[\hat{\theta}(\mathbf{X})] \succcurlyeq \tau'(\theta) I_{\mathbf{X}}^{-1}(\theta) (\tau'(\theta))^\top.$$

Доказательство. Рассуждая ровно так же, как и в одномерном случае, получаем, что $\mathbb{E}_\theta[U_\theta(\mathbf{X})] = \mathbf{0}$ и

$$\tau'(\theta) = \mathbb{E}_\theta[\hat{\theta}(\mathbf{X})(U_\theta(\mathbf{X}))^\top] = \mathbb{E}_\theta[(\hat{\theta}(\mathbf{X}) - \tau(\theta))(U_\theta(\mathbf{X}))^\top].$$

Далее, по матричному неравенству Коши-Буняковского-Шварца получаем, что

$$\mathbf{D}_\theta[\hat{\theta}(\mathbf{X})] \succcurlyeq \tau'(\theta) I_{\mathbf{X}}^{-1}(\theta) (\tau'(\theta))^\top.$$

□

3 Лекция 3

3.1 Оценки максимального правдоподобия

Сегодня мы будем обсуждать то, что преобразило статистику — метод максимального правдоподобия. Начнём с того, что введём понятие правдоподобия.

Определение 17. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — наблюдение с неизвестным распределением $\mathbf{P} \in \{\mathbf{P}_\theta, \theta \in \Theta\}$, где $\{\mathbf{P}_\theta, \theta \in \Theta\}$ есть доминируемое семейство с плотностью $p_\theta(\mathbf{x})$.⁶ Тогда функцией правдоподобия называется случайная величина $f_\theta(\mathbf{X}) = p_\theta(\mathbf{X})$.

Примечание. Если $\mathbf{X} = (X_1, \dots, X_n)$ — выборка, то плотность случайного вектора разбивается в произведение плотностей координат:

$$f_\theta(\mathbf{X}) = p_\theta(\mathbf{X}) = \prod_{i=1}^n p_\theta(X_i).$$

Теперь можно ввести и сам метод.

Определение 18. Оценкой параметра θ по методу максимального правдоподобия, или же *оценкой максимального правдоподобия*, называется

$$\hat{\theta}(\mathbf{X}) = \arg \max_{\theta \in \Theta} f_\theta(\mathbf{X}).$$

⁶Плотность понимается в обобщённом смысле — за подробностями обращайтесь ко второй лекции.

Данное определение уже накладывает несколько ограничений: как минимум, то, что максимум существует и он единственен. Но в дальнейшем будем считать, что они выполнены. Философия этого метода уже не так очевидна. Она состоит в том, что «мы живём в наиболее вероятном мире».

Допустим, что мы рассматриваем схему Бернулли и у нас выпало много нулей и мало единиц. Мы думаем: «Наверное, это неспроста!» Тогда, наверное, так и должно быть на самом деле — действительно, в дискретном случае функция правдоподобия есть вероятность того, что выпадает данный набор. Далее, мы подбираем θ такое, что данный набор наблюдений наиболее вероятен — коли уж он выпал, то истинным значением должно быть только то, в котором он выпадает с наибольшей вероятностью.

Рассмотрим пару примеров нахождения оценки максимального правдоподобия, для того, чтобы понять, что вообще происходит.

Задача 4. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из равномерного распределения $U(0, \theta)$, $\theta > 0$. Найти оценку максимального правдоподобия параметра θ .

Решение. Начнём с того, что распишем функцию правдоподобия:

$$f_{\theta}(\mathbf{X}) = \prod_{i=1}^n p_{\theta}(X_i) = \prod_{i=1}^n \frac{1}{\theta} [0 \leq X_i \leq \theta] = \frac{1}{\theta^n} [0 \leq X_{(1)} \leq X_{(n)} \leq \theta].$$

Теперь нам нужно максимизировать её, как функцию от θ . Заметим, что θ^{-n} есть монотонно убывающая функция, поэтому нужно взять минимальное θ такое, что функция правдоподобия не обратится в ноль. Но тогда $\hat{\theta}(\mathbf{X}) = X_{(n)}$. \square

Задача 5. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из нормального распределения $\mathcal{N}(\mu, \sigma^2)$. Найдите оценку максимального правдоподобия параметра $\theta = (\mu, \sigma^2)$.

Решение. Опять же, начнём с того, что распишем функцию правдоподобия:

$$f_{\theta}(\mathbf{X}) = \prod_{i=1}^n p_{\theta}(X_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(X_i - \mu)^2}{2\sigma^2} \right\} = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right\}.$$

Нам нужно максимизировать её одновременно по μ и по σ . Но работать с экспонентой достаточно неудобно, поэтому прологарифмируем её. Это не поменяет решения, так как логарифм биективно переводит \mathbb{R}_{++} в \mathbb{R} :

$$\ln f_{\theta}(\mathbf{X}) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2.$$

Эта функция дифференцируема, поэтому можно достаточно легко найти всех претендентов на точки экстремума приравняв производной к нулю:

$$\begin{aligned} 0 &= \frac{\partial}{\partial \mu} \ln f_{\theta}(\mathbf{X}) = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu), \\ 0 &= \frac{\partial}{\partial \sigma^2} \ln f_{\theta}(\mathbf{X}) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (X_i - \mu)^2. \end{aligned}$$

Отсюда несложно получить, что претендентами на точки экстремума будут

$$\begin{aligned} \hat{\mu}(\mathbf{X}) &= \frac{1}{n} \sum_{i=1}^n X_i = \bar{\mathbf{X}}, \\ \hat{\sigma}^2(\mathbf{X}) &= \frac{1}{n} \sum_{i=1}^n \left(X_i - \frac{1}{n} \sum_{i=1}^n X_i \right)^2 = S^2. \end{aligned}$$

Но данная точка действительно будет являться максимумом, что несложно проверить по значениям производных. Так как она есть единственный максимум, то это глобальный максимум и $\hat{\theta}(\mathbf{X}) = (\bar{\mathbf{X}}, S^2)$. \square

Примечание. Несложно показать, что S^2 — смещённая оценка. Но оценка максимального правдоподобия не обязательно несмещённая. Вообще, несмещённость — достаточно сильное свойство (и получить из оценки несмещённую может быть весьма нетривиально) в том плане, что мы хотим равенство матожиданий. Но у ОМП есть так называемая *асимптотическая несмещённость*: предел матожиданий будет именно таким, каким и должен быть.

У оценок максимального правдоподобия и функции правдоподобия есть достаточно интересные свойства. Но они требуют некоторых условий регулярности. Будем постепенно формулировать их и доказывать свойства.

(R1) Параметрическое семейство распределений $\{P_\theta \mid \theta \in \Theta\}$ — это доминируемое семейство с плотностью $p_\theta(x)$ и *различимыми распределениями*, то есть $P_{\theta_0} = P_{\theta_1}$ почти везде тогда и только тогда, когда $\theta_0 = \theta_1$.

(R2) $\mathbf{X} = (X_1, \dots, X_n)$ — выборка растущего размера из неизвестного распределения $P \in \{P_\theta \mid \theta \in \Theta\}$.

(R3) $A = \{x: p_\theta(x) > 0\}$ не зависит от θ .

Теорема 6 (Экстремальное свойство правдоподобия). *В условиях регулярности (R1)–(R3) для всех различных $\theta_0, \theta_1 \in \Theta$*

$$\lim_{n \rightarrow \infty} P_{\theta_0}(f_{\theta_0}(\mathbf{X}) > f_{\theta_1}(\mathbf{X})) = 1.$$

Доказательство. Будем считать, что мы будем работать в A . Посмотрим, при каких условиях выполняется событие $f_{\theta_0}(\mathbf{X}) > f_{\theta_1}(\mathbf{X})$. Для этого прологарифмируем и преобразуем выражение:

$$\ln \frac{f_{\theta_1}(\mathbf{X})}{f_{\theta_0}(\mathbf{X})} < 0 \implies \frac{1}{n} \sum_{i=1}^n \ln \frac{p_{\theta_1}(X_i)}{p_{\theta_0}(X_i)} < 0.$$

По усиленному закону больших чисел

$$\frac{1}{n} \sum_{i=1}^n \ln \frac{p_{\theta_1}(X_i)}{p_{\theta_0}(X_i)} \xrightarrow{P_{\theta_0}\text{-п.н.}} E_{\theta_0} \left[\ln \frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} \right].$$

Теперь докажем, что⁷

$$E_{\theta_0} \left[\ln \frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} \right] < 0.$$

Воспользуемся неравенством Йенсена:

$$E_{\theta_0} \left[\ln \frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} \right] \leq \ln E_{\theta_0} \left[\frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} \right] = \ln \int_A \frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} p_{\theta_0}(X_1) \mu(dx) = \ln E_{\theta_1}[1] = 0.$$

Но почему оно не равно нулю? Предположим, что это так:

$$E_{\theta_0} \left[\ln \frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} \right] = 0.$$

Но в таком случае можно воспользоваться критерием равенства для неравенства Йенсена: $\phi(E[\xi]) = E[\phi(\xi)]$ тогда и только тогда, когда ϕ линейна почти везде. Но $\ln(x)$ нелинейна. Тогда получаем, что аргумент должен быть равен единице почти везде: $\mu(\{x: p_{\theta_0}(x) = p_{\theta_1}(x)\}) = 1$. Но это означает, что $\theta_0 = \theta_1$, что противоречит условию.

В итоге получаем, что

$$\lim_{n \rightarrow \infty} P_{\theta_0}(f_{\theta_0}(\mathbf{X}) > f_{\theta_1}(\mathbf{X})) = P_{\theta_0} \left(E_{\theta_0} \left[\ln \frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} \right] < 0 \right) = 1. \quad \square$$

Следствие (Состоятельность оценки максимального правдоподобия). Если Θ конечно, то оценка максимального правдоподобия состоятельна.

Доказательство. Пусть $\hat{\theta}_n(\mathbf{X})$ — это оценка максимального правдоподобия. Тогда по экстремальному свойству правдоподобия для любого $\theta_0 \in \Theta$

$$\lim_{n \rightarrow \infty} P_{\theta_0}(\hat{\theta}_n(\mathbf{X}) = \theta_0) = \lim_{n \rightarrow \infty} P_{\theta_0}(\forall \theta \neq \theta_0, f_{\theta_0}(\mathbf{X}) > f_{\theta}(\mathbf{X})) = 1. \quad \square$$

Введём ещё два условия регулярности:

(R4) Θ есть открытый интервал на \mathbb{R} .

⁷Небольшое примечание: если добавить минус к этому матожиданию, то получим широко известную *дивергенцию Кульбака-Лейблера*. По сути, мы доказываем то, что она неотрицательна и то, что она равна нулю тогда и только тогда, когда плотности равны почти везде.

(R5) $p_\theta(x)$ непрерывно дифференцируема по θ для всех $x \in A$.

Теперь можно доказать хорошее свойство, плотно связанное с оценками максимального правдоподобия.

Теорема 7 (Состоятельность решения уравнения правдоподобия). *В условиях регулярности (R1)–(R5) уравнение правдоподобия*

$$\frac{\partial}{\partial \theta} \ln f_\theta(\mathbf{X}) = 0$$

с вероятностью, стремящейся к 1, имеет решение, которое сходится по вероятности к истинному значению параметра.

Доказательство. Пусть θ_0 есть истинное значение параметра. Возьмём $\delta > 0$ такое, что $[\theta_0 - \delta, \theta_0 + \delta] \subset \Theta$ (это возможно из-за открытости Θ). Далее, введём следующее событие:

$$A_n = \{ f_{\theta_0}(\mathbf{X}) > f_{\theta_0 + \delta}(\mathbf{X}), f_{\theta_0}(\mathbf{X}) > f_{\theta_0 - \delta}(\mathbf{X}) \}$$

Тогда по экстремальному свойству правдоподобия

$$\lim_{n \rightarrow \infty} P_{\theta_0}(A_n) = 1.$$

Что можно сказать, если выполнено A_n ? Так как по (R6) производная логарифма функции правдоподобия непрерывна по θ , то на отрезке $[\theta_0 - \delta, \theta_0 + \delta]$ будет точка, в которой возрастание заменяется убыванием. Следовательно, на нём будет хотя бы один корень уравнения правдоподобия. Допустим, что на этом отрезке есть несколько корней (не обязательно конечное число). Пусть $\tilde{\theta}(\mathbf{X})$ — ближайший к θ_0 корень. Это возможно, так как предел корней тоже является корнем (так как производная непрерывна). Тогда оказывается, что для всех $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P_{\theta_0}(|\tilde{\theta}(\mathbf{X}) - \theta_0| \leq \varepsilon) = 1.$$

Почему это так? Зафиксируем ε и заметим, что рассуждения выше легальны для $\delta = \varepsilon$. Следовательно, с вероятностью, стремящейся к 1, на отрезке $[\theta_0 - \varepsilon, \theta_0 + \varepsilon]$ будет корень. Но $\tilde{\theta}(\mathbf{X})$ — ближайший к θ_0 корень. Тогда он лежит в этом отрезке. \square

Вот доказали мы эту теорему. Но она больше напоминает решето, из которого вытекают проблемы. Почему?

1. Например, корней уравнения правдоподобия может быть несколько. В доказательстве мы выбираем ближайший из них к истинному значению. Но как его выбрать?
2. Даже если мы его найдём, то он зависит от истинного значения, то есть вообще не является оценкой.
3. Почему $\tilde{\theta}(\mathbf{X})$ есть точка максимума? Мы сказали, что это корень уравнения, но он не обязательно даёт максимум — может оказаться, что это точка минимума или же точка перегиба (если корней несколько).
4. Корень существует не всегда, а только с большой вероятностью.

Впрочем, если сказать, что уравнение правдоподобия имеет только один корень, то всё достаточно неплохо. Четвёртый и второй вопросы отпадают сразу же, да и первый тоже. Остался третий, но он тоже исправляется:

Теорема 8. *Если для всех $\mathbf{X} = (X_1, \dots, X_n)$ уравнение правдоподобия имеет единственный корень $\hat{\theta}(\mathbf{X})$, то с вероятностью, стремящейся к 1, $\hat{\theta}(\mathbf{X})$ будет оценкой максимального правдоподобия и оценка максимального правдоподобия будет состоятельной.*

Доказательство. По сути, доказательство повторяет то, что было сказано выше. Пусть θ_0 — истинное значение параметра. Снова возьмём $\delta > 0$ такое, что $[\theta_0 - \delta, \theta_0 + \delta] \subset \Theta$ и заметим, что если

$$\lim_{n \rightarrow \infty} P_{\theta_0}(A_n) = 1, \text{ где } A_n = \{ f_{\theta_0}(\mathbf{X}) > f_{\theta_0 + \delta}(\mathbf{X}), f_{\theta_0}(\mathbf{X}) > f_{\theta_0 - \delta}(\mathbf{X}) \}.$$

Однако если выполнено A_n , то внутри $[\theta_0 - \delta, \theta_0 + \delta]$ есть точка локального максимума. Как известно, в ней производная равна нулю, и, следовательно, она будет корнем уравнения правдоподобия. Но тогда эта точка есть $\hat{\theta}(\mathbf{X})$. Далее, это должен быть глобальный максимум, так как иначе найдётся локальный минимум, что противоречит единственности корня. Следовательно, если выполнено A_n , то $\hat{\theta}(\mathbf{X})$ есть оценка максимального правдоподобия. Тогда

$$\lim_{n \rightarrow \infty} P_{\theta_0}(\hat{\theta}(\mathbf{X}) = \text{ОМП}) = 1.$$

Но, как известно, $\hat{\theta}(\mathbf{X})$ есть ближайший к θ_0 корень. Тогда для всех $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta_0}(|\hat{\theta}(\mathbf{X}) - \theta_0| \leq \varepsilon) = 1.$$

Тем самым получаем, что и ОМП будет состоятельной оценкой параметра θ . \square

Поехали дальше. Эти условия регулярности уже далеко не так очевидны на первый взгляд.

(R6) $p_\theta(x)$ трижды непрерывно дифференцируема по θ для всех $x \in A$.

(R7) Интеграл

$$\int_A p_\theta(x) \mu(dx)$$

можно дважды дифференцировать под знаком интеграла.

(R8) Для всех $\theta \in \Theta$

$$0 < i(\theta) = \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \ln p_\theta(X_1) \right)^2 \right] < +\infty.$$

(R9) Для любого $\theta_0 \in \Theta$ существует $\delta > 0$ и функция $M(x)$ такая, что для всех $\theta \in [\theta_0 - \delta, \theta_0 + \delta]$

$$\left| \frac{\partial^3}{\partial \theta^3} \ln p_\theta(x) \right| \leq M(x), \text{ причём } \mathbb{E}_{\theta_0}[M(X_1)] < +\infty.$$

Теорема 9. В условиях регулярности (R1)–(R9) любая состоятельная последовательность $\{\hat{\theta}_n(\mathbf{X}) \mid n \in \mathbb{N}\}$ корней уравнения правдоподобия удовлетворяет свойству асимптотической нормальности: для всех $\theta_0 \in \Theta$

$$\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0) \xrightarrow{d_{\theta_0}} \mathcal{N}\left(0, \frac{1}{i(\theta_0)}\right).$$

Доказательство. Обозначим $\mathcal{L}(\mathbf{X}, \theta) = \ln f_\theta(\mathbf{X})$ и будем говорить, что $\mathcal{L}^{(n)}(\mathbf{X}, \theta)$ есть n -я частная производная $\mathcal{L}(\mathbf{X}, \theta)$ по θ . Далее, пусть θ_0 есть истинное значение параметра, то есть $\hat{\theta}_n(\mathbf{X})$ сходится к θ_0 по вероятности \mathbb{P}_{θ_0} . Разложим $\mathcal{L}'(\mathbf{X}, \theta)$ в ряд Тейлора в точке θ_0 :

$$\mathcal{L}'(\mathbf{X}, \theta) = \mathcal{L}'(\mathbf{X}, \theta_0) + \mathcal{L}''(\mathbf{X}, \theta_0)(\theta - \theta_0) + \frac{1}{2} \mathcal{L}'''(\mathbf{X}, \tilde{\theta})(\theta - \theta_0)^2,$$

где $\tilde{\theta}$ находится между θ и θ_0 . Теперь подставим $\theta = \hat{\theta}_n(\mathbf{X})$:

$$\mathcal{L}'(\mathbf{X}, \hat{\theta}_n(\mathbf{X})) = \mathcal{L}'(\mathbf{X}, \theta_0) + \mathcal{L}''(\mathbf{X}, \theta_0)(\hat{\theta}_n(\mathbf{X}) - \theta_0) + \frac{1}{2} \mathcal{L}'''(\mathbf{X}, \tilde{\theta}_n)(\hat{\theta}_n(\mathbf{X}) - \theta_0)^2,$$

где $\tilde{\theta}_n$ находится между $\hat{\theta}_n(\mathbf{X})$ и θ_0 . Теперь вспомним, что $\hat{\theta}_n(\mathbf{X})$ есть корень уравнения правдоподобия. Тогда $\mathcal{L}'(\mathbf{X}, \hat{\theta}_n(\mathbf{X})) = 0$. Теперь выразим $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$ из данного равенства следующим образом:

$$\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0) = \frac{-\sqrt{n}\mathcal{L}'(\mathbf{X}, \theta_0)}{\mathcal{L}''(\mathbf{X}, \theta_0) + \frac{1}{2}\mathcal{L}'''(\mathbf{X}, \tilde{\theta}_n)(\hat{\theta}_n(\mathbf{X}) - \theta_0)} = \frac{-\frac{1}{\sqrt{n}}\mathcal{L}'(\mathbf{X}, \theta_0)}{\frac{1}{n}\mathcal{L}''(\mathbf{X}, \theta_0) + \frac{1}{2n}\mathcal{L}'''(\mathbf{X}, \tilde{\theta}_n)(\hat{\theta}_n(\mathbf{X}) - \theta_0)}$$

Теперь рассмотрим это выражение по частям.

- Начнём с числителя. Заметим, что

$$-\frac{1}{\sqrt{n}}\mathcal{L}'(\mathbf{X}, \theta_0) = -\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln p_\theta(X_i) \Big|_{\theta=\theta_0} = -\sqrt{n} \sum_{i=1}^n U_{\theta_0}(X_i)$$

Как известно, возможность дифференцирования под знаком интеграла показывает, что $\mathbb{E}_{\theta_0}[U_{\theta_0}(X_1)] = 0$. Тогда по центральной предельной теореме

$$-\frac{1}{\sqrt{n}} \sum_{i=1}^n U_{\theta_0}(X_i) \xrightarrow{d_{\theta_0}} -\mathcal{N}(0, \mathbb{D}_{\theta_0}[U_{\theta_0}(X_1)]) = -\mathcal{N}(0, i(\theta_0)).$$

- По усиленному закону больших чисел

$$\frac{1}{n} \mathcal{L}''(\mathbf{X}, \theta_0) \xrightarrow{\text{P}_{\theta_0}\text{-п.н.}} \mathbb{E}_{\theta_0} \left[\left. \frac{\partial^2}{\partial \theta^2} \ln p_{\theta}(X_1) \right|_{\theta=\theta_0} \right].$$

Докажем, что

$$-\mathbb{E}_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \ln p_{\theta}(X_1) \right] = i(\theta).$$

Доказательство. Заметим, что

$$\begin{aligned} \frac{\partial^2}{\partial \theta^2} \ln p_{\theta}(x) &= \frac{\partial}{\partial \theta} \left(\frac{1}{p_{\theta}(x)} \frac{\partial p_{\theta}(x)}{\partial \theta} \right) = -\frac{1}{(p_{\theta}(x))^2} \left(\frac{\partial p_{\theta}(x)}{\partial \theta} \right)^2 + \frac{1}{p_{\theta}(x)} \frac{\partial^2 p_{\theta}(x)}{\partial \theta^2} \\ &= -\left(\frac{\partial}{\partial \theta} \ln p_{\theta}(x) \right)^2 + \frac{1}{p_{\theta}(x)} \frac{\partial^2 p_{\theta}(x)}{\partial \theta^2}. \end{aligned}$$

Теперь возьмём матожидание и воспользуемся тем, что мы можем два раза дифференцировать под знаком интеграла:

$$\begin{aligned} \mathbb{E}_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \ln p_{\theta}(X_1) \right] &= \int_A \frac{\partial^2 \ln p_{\theta}(x)}{\partial \theta^2} p_{\theta}(x) \mu(dx) \\ &= - \int_A \left(\frac{\partial \ln p_{\theta}(x)}{\partial \theta} \right)^2 p_{\theta}(x) \mu(dx) + \int_A \frac{\partial^2 p_{\theta}(x)}{\partial \theta^2} \mu(dx) \\ &= -i(\theta) + \frac{\partial^2}{\partial \theta^2} \int_A p_{\theta}(x) \mu(dx) = -i(\theta) \end{aligned} \quad \square$$

Отсюда получаем, что

$$\frac{1}{n} \mathcal{L}''(\mathbf{X}, \theta_0) \xrightarrow{\text{P}_{\theta_0}\text{-п.н.}} -i(\theta_0).$$

- Теперь перейдём к третьему члену и покажем, что он стремится к нулю по вероятности. Заметим, что по условию $\hat{\theta}(\mathbf{X})$ сходится к θ_0 по вероятности. Из этого можно сделать вывод, что $\tilde{\theta}_n$ тоже стремится к θ_0 по вероятности. Далее, по (R9)

$$\left| \frac{1}{n} \mathcal{L}'''(\mathbf{X}, \theta) \right| \leq \frac{1}{n} \sum_{i=1}^n M(X_i),$$

что имеет конечное матожидание. Тогда получаем произведение ограниченной случайной величины на нечто, что сходится к нулю по вероятности. Это сходится к нулю по вероятности.

Комбинируя вышесказанное, по лемме Слущкого получаем, что

$$\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0) \xrightarrow{d_{\theta_0}} \frac{1}{i(\theta_0)} \mathcal{N}(0, i(\theta_0)) = \mathcal{N}\left(0, \frac{1}{i(\theta_0)}\right). \quad \square$$

И получаем приятное с практической точки зрения свойство:

Следствие. Если в условиях теоремы для всех $\mathbf{X} = (X_1, \dots, X_n)$ существует единственное решение уравнения правдоподобия, то оно является оценкой максимального правдоподобия, причём ОМП будет асимптотически нормальной оценкой параметра θ с асимптотической дисперсией $i^{-1}(\theta)$.

Оказывается, что можно предложить нижнюю границу не только для обычной дисперсии (что даёт неравенство Рао-Крамера), но и для асимптотической дисперсии. Этот результат называется *теоремой Бахадура*. Сформулируем его:

Теорема 10 (Бахадур). Если в условиях регулярности (R1)–(R9) оценка $\hat{\theta}(\mathbf{X})$ является асимптотически нормальной оценкой параметра θ с асимптотической дисперсией $\sigma^2(\theta)$, то $\sigma^2(\theta) \geq i^{-1}(\theta)$ почти везде (то есть неравенство нарушается только на множестве лебеговой меры 0).

Ну и сразу же пример, в котором ограничение выполнено не везде, но почти везде.

Задача 6. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из нормального распределения $\mathcal{N}(\theta, 1)$. Введём следующую оценку параметра θ :

$$\hat{\theta}_n(\mathbf{X}) = \begin{cases} \bar{X}, & |\bar{X}| \geq n^{-1/4} \\ \bar{X}/2, & |\bar{X}| < n^{-1/4} \end{cases}$$

Найдите асимптотическую дисперсию $\sigma^2(\theta)$ и сравните её с обратной информацией Фишера $i^{-1}(\theta)$ одного элемента.

Решение. Как известно, усиленный закон больших чисел имеет скорость сходимости порядка $\mathcal{O}(n^{-1/2})$ с вероятностью, стремящейся к 1. Тогда если $\theta \neq 0$, то

$$\lim_{n \rightarrow \infty} P_\theta(\hat{\theta}_n(\mathbf{X}) = \bar{X}) = 1$$

Теперь рассмотрим распределение $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta)$. Для этого воспользуемся формулой полной вероятности:

$$\begin{aligned} P(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x) &= P(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x \mid \hat{\theta}_n(\mathbf{X}) = \bar{X})P(\hat{\theta}_n(\mathbf{X}) = \bar{X}) \\ &\quad + P(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x \mid \hat{\theta}_n(\mathbf{X}) = \bar{X}/2)P(\hat{\theta}_n(\mathbf{X}) = \bar{X}/2) \\ &= P(\sqrt{n}(\bar{X} - \theta) \leq x)P(\hat{\theta}_n(\mathbf{X}) = \bar{X}) \\ &\quad + P(\sqrt{n}(\bar{X}/2 - \theta) \leq x)P(\hat{\theta}_n(\mathbf{X}) = \bar{X}/2) \end{aligned}$$

Теперь устремим n к бесконечности. Согласно результату выше и центральной предельной теореме получаем, что

$$\lim_{n \rightarrow \infty} P(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x) = P(\mathcal{N}(0, 1) \leq x) = \Phi(x) \implies \sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \xrightarrow{d_\theta} \mathcal{N}(0, 1).$$

Далее скажем, что пусть $\theta = 0$. В таком случае утверждение о скорости сходимости усиленного закона больших чисел говорит, что

$$\lim_{n \rightarrow \infty} P_\theta(\hat{\theta}_n(\mathbf{X}) = \bar{X}/2) = 1.$$

Применяя формулу полной вероятности, получаем, что

$$P(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x) = P(\sqrt{n}\bar{X} \leq x)P(\hat{\theta}_n(\mathbf{X}) = \bar{X}) + P(\sqrt{n}\bar{X} \leq 2x)P(\hat{\theta}_n(\mathbf{X}) = \bar{X}/2).$$

Следовательно,

$$\lim_{n \rightarrow \infty} P(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x) = \Phi(2x) \implies \sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \xrightarrow{d_\theta} \mathcal{N}(0, 1/4).$$

Однако информация Фишера равна

$$i(\theta) = E_\theta[(X_1 - \theta)^2] = D_\theta[X_1] = 1.$$

Тогда мы получаем, что везде, кроме нуля, $\sigma^2(\theta) \geq i^{-1}(\theta)$. □

4 Лекция 4

4.1 Теорема Бахадура и следствия из неё

На прошлой лекции мы сформулировали теорему Бахадура. Сейчас мы попытаемся доказать её. Но для начала потребуем выполнение следующих условий регулярности:

- (R1) $\mathbf{X} = (X_1, \dots, X_n)$ — выборка растущего размера из неизвестного распределения $P \in \{P_\theta, \theta \in \Theta\}$, где $\{P_\theta, \theta \in \Theta\}$ — доминируемое семейство с плотностью $p_\theta(x)$ по мере μ .
- (R2) Будем считать, что θ есть одномерный параметр и Θ есть открытый интервал на \mathbb{R} .
- (R3) Множество $A = \{x: p_\theta(x) > 0\}$ не зависит от θ .
- (R4) Функция $\ell(\theta, x) = \ln p_\theta(x)$ дважды непрерывно дифференцируема по θ для всех $x \in A$. Для удобства будем считать, что

$$\ell^{(n)}(\theta, x) = \frac{\partial^n}{\partial \theta^n} \ell(\theta, x).$$

(R5) Для всех $\theta \in \Theta$ $E_\theta[\ell'(\theta, X_1)] = 0$ и $0 < E_\theta[(\ell'(\theta, X_1))^2] = i(\theta) < +\infty$. Стоит заметить, что $-E_\theta[\ell''(\theta, X_1)] = i(\theta)$.

(R6) Для всех $\theta_0 \in \Theta$ существует $\delta > 0$ и функция $M(x)$ такая, что для всех $\theta \in [\theta_0 - \delta, \theta_0 + \delta]$ $|\ell'(\theta, x)| \leq M(x)$ и $E_\theta[M(X_1)] < +\infty$.

Теорема 11. Пусть в условиях регулярности (R1)–(R6) для последовательности оценок $t_n(\mathbf{X})$ параметра θ выполнено условие асимптотической нормальности: для всех $\theta \in \Theta$

$$\sqrt{n}(t_n(\mathbf{X}) - \theta) \xrightarrow{d_\theta} \mathcal{N}(0, \sigma^2(\theta)).$$

Далее, зафиксируем $\theta_0 \in \Theta$ и положим $\theta_n = \theta_0 + n^{-1/2}$. Если

$$\lim_{n \rightarrow \infty} P_{\theta_n}(t_n(\mathbf{X}) \leq \theta_n) \leq \frac{1}{2},$$

то $\sigma^2(\theta_0) \geq i^{-1}(\theta)$.

Для доказательства данной теоремы нужно сперва доказать несколько вспомогательных утверждений и ввести несколько обозначений. Начнём с того, что введём логарифмическую функцию правдоподобия:

$$\mathcal{L}(\theta, \mathbf{X}) = \ln f_\theta(\mathbf{X}) = \sum_{i=1}^n \ell(\theta, X_i).$$

Далее введём следующую случайную величину:

$$T_n = \frac{1}{\sqrt{i(\theta_0)}} \left(\mathcal{L}(\theta_n, \mathbf{X}) - \mathcal{L}(\theta_0, \mathbf{X}) + \frac{i(\theta_0)}{2} \right).$$

Теперь покажем, что для случайной величины T_n выполнено одно интересное свойство:

Лемма. При фиксированном θ_0 последовательность случайных величин T_n сходится по распределению к стандартному нормальному распределению: $T_n \xrightarrow{d_{\theta_0}} \mathcal{N}(0, 1)$.

Доказательство. Воспользуемся дважды непрерывной дифференцируемостью $\ell(\theta, x)$ и разложим $\mathcal{L}(\theta_n, \mathbf{X})$ по θ в окрестности θ_0 :

$$\begin{aligned} \mathcal{L}(\theta_n, \mathbf{X}) &= \mathcal{L}(\theta_0, \mathbf{X}) + \mathcal{L}'(\theta_0, \mathbf{X})(\theta_n - \theta_0) + \frac{1}{2} \mathcal{L}''(\tilde{\theta}_n, \mathbf{X})(\theta_n - \theta_0)^2 \\ &= \mathcal{L}(\theta_0, \mathbf{X}) + \frac{1}{\sqrt{n}} \mathcal{L}'(\theta_0, \mathbf{X}) + \frac{1}{2n} \mathcal{L}''(\tilde{\theta}_n, \mathbf{X}), \end{aligned}$$

где $\tilde{\theta}_n \in (\theta_0, \theta_n)$. Далее, введём следующую случайную величину:

$$\xi_n = \frac{1}{2n} \left(\mathcal{L}''(\tilde{\theta}_n, \mathbf{X}) - \mathcal{L}''(\theta_0, \mathbf{X}) \right).$$

Докажем, что эта последовательность стремится к нулю почти наверное: $\xi_n \xrightarrow{P_{\theta_0}\text{-п.н.}} 0$. Для этого введём следующую функцию:

$$A(x, \delta) = \sup_{\theta \in [\theta_0 - \delta, \theta_0 + \delta]} |\ell''(\theta, x) - \ell''(\theta_0, x)|.$$

Далее, положим $m(\delta) = E_{\theta_0}[A(X_1, \delta)]$. Заметим, что если $\delta_1 > \delta_2$, то $A(x, \delta_1) \geq A(x, \delta_2)$ и при $\delta \rightarrow 0$ $A(x, \delta) \rightarrow 0$. Тогда $A(x, \delta)$ монотонно убывает к нулю при $\delta \rightarrow 0$. Далее, воспользуемся свойством регулярности (R6):

$$A(x, \delta) = \sup_{\theta \in [\theta_0 - \delta, \theta_0 + \delta]} |\ell''(\theta, x) - \ell''(\theta_0, x)| \leq \sup_{\theta \in [\theta_0 - \delta, \theta_0 + \delta]} 2M(x) = 2M(x).$$

Тем самым можно применить теорему Лебега о мажорируемой сходимости. Из неё следует, что $m(\delta) \rightarrow 0$ при $\delta \rightarrow 0$. Далее, заметим, что если взять $\delta > n^{-1/2}$, то

$$|\xi_n| = \frac{1}{2n} \left| \sum_{i=1}^n (\ell''(\tilde{\theta}_n, X_i) - \ell''(\theta_0, X_i)) \right| \leq \frac{1}{2n} \sum_{i=1}^n |\ell''(\tilde{\theta}_n, X_i) - \ell''(\theta_0, X_i)| \leq \frac{1}{2n} \sum_{i=1}^n A(X_i, \delta).$$

Однако по усиленному закону больших чисел

$$\frac{1}{2n} \sum_{i=1}^n A(X_i, \delta) \xrightarrow{P_{\theta_0}\text{-п.н.}} \frac{1}{2} m(\delta).$$

Из вышесказанного следует, что для любого $\delta > 0$

$$\overline{\lim}_{n \rightarrow \infty} |\xi_n| \leq \frac{1}{2} m(\delta) \quad P_{\theta_0}\text{-п.н.}$$

Следовательно, $\xi_n \xrightarrow{P_{\theta_0}\text{-п.н.}} 0$.

Теперь преобразуем T_n следующим образом:

$$\begin{aligned} T_n &= \frac{1}{\sqrt{i(\theta_0)}} \left(\mathcal{L}(\theta_n, \mathbf{X}) - \mathcal{L}(\theta_0, \mathbf{X}) + \frac{i(\theta_0)}{2} \right) \\ &= \frac{1}{\sqrt{i(\theta_0)}} \left(\mathcal{L}(\theta_0, \mathbf{X}) + \frac{1}{\sqrt{n}} \mathcal{L}'(\theta_0, \mathbf{X}) + \frac{1}{2n} \mathcal{L}''(\tilde{\theta}_n, \mathbf{X}) - \mathcal{L}(\theta_0, \mathbf{X}) + \frac{i(\theta_0)}{2} \right) \\ &= \frac{1}{\sqrt{i(\theta_0)}} \left(\frac{1}{\sqrt{n}} \mathcal{L}'(\theta_0, \mathbf{X}) + \xi_n + \frac{1}{2n} \mathcal{L}''(\theta_0, \mathbf{X}) + \frac{i(\theta_0)}{2} \right). \end{aligned}$$

Осталось показать, что это стремится по распределению туда, куда нужно.

- Начнём с члена $n^{-1/2} \mathcal{L}'(\theta_0, \mathbf{X})$. Заметим, что из условия регулярности (R5) и центральной предельной теоремы следует, что

$$\frac{1}{\sqrt{n}} \mathcal{L}'(\theta_0, \mathbf{X}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \ell'(\theta_0, X_i) \xrightarrow{d_{\theta_0}} \mathcal{N}(0, i(\theta_0)).$$

- Теперь рассмотрим член $(2n)^{-1} \mathcal{L}''(\theta_0, \mathbf{X})$. По усиленному закону больших чисел

$$\frac{1}{2n} \mathcal{L}''(\theta_0, \mathbf{X}) = \frac{1}{2n} \sum_{i=1}^n \ell''(\theta_0, X_i) \xrightarrow{P_{\theta_0}\text{-п.н.}} \frac{1}{2} E_{\theta_0}[\ell''(\theta_0, X_1)] = -\frac{i(\theta_0)}{2}.$$

- Из пунктов выше следует, что

$$\xi_n + \frac{1}{2n} \mathcal{L}''(\theta_0, \mathbf{X}) + \frac{i(\theta_0)}{2} \xrightarrow{P_{\theta_0}\text{-п.н.}} 0.$$

В итоге по лемме Slutsky получаем, что

$$T_n \xrightarrow{d_{\theta_0}} \frac{1}{\sqrt{i(\theta_0)}} \mathcal{N}(0, i(\theta_0)) = \mathcal{N}(0, 1). \quad \square$$

Для доказательства следующей леммы нам понадобится достаточно известный факт, связанный с равномерной интегрируемостью.⁸ Но перед этим докажем одно простое утверждение:

Лемма (Абсолютная непрерывность математического ожидания). Пусть ξ — интегрируемая случайная величина. Тогда для любого $\varepsilon > 0$ существует $\delta > 0$ такое, что для любого события F такого, что $P(F) < \delta$, выполнено $E[|\xi| \mathbf{1}_F] < \varepsilon$.

Доказательство. Пусть это не так. Тогда существует $\varepsilon > 0$ и последовательность событий F_n такая, что $P(F_n) < 2^{-n}$ и $E[|\xi| \mathbf{1}_{F_n}] \geq \varepsilon$. Но $|\xi| \mathbf{1}_{F_n} \rightarrow 0$ почти наверное и $|\xi| \mathbf{1}_{F_n} \leq |\xi|$. Следовательно, по теореме Лебега о мажорируемой сходимости $\varepsilon \leq E[|\xi| \mathbf{1}_{F_n}] \rightarrow 0$. Противоречие. \square

Теорема 12. Пусть ξ_n — последовательность неотрицательных интегрируемых случайных величин таких, что ξ_n сходится к ξ по распределению. Тогда $E[\xi_n] \rightarrow E[\xi]$ тогда и только тогда, когда случайные величины ξ_n равномерно интегрируемы, то есть

$$\lim_{c \rightarrow +\infty} \sup_{n \in \mathbb{N}} E[\xi_n \mathbf{1}_{\xi_n \geq c}] = 0.$$

⁸Доказательство этого факта — дополнительный материал, который на лекции был оставлен в качестве упражнения.

Доказательство. Начнём с доказательства в прямую сторону: то есть докажем, что из сходимости матожиданий следует равномерная интегрируемость. Заметим, что $E[|\xi_n - \xi|] \rightarrow 0$, так как случайные величины неотрицательны (и, следовательно, сходятся в L_1). Теперь нам нужно показать, что для любого $\varepsilon > 0$ можно подобрать $c > 0$ такое, что для всех $n \in \mathbb{N}$ $E[\xi_n[\xi_n \geq c]] < 2\varepsilon$.

Хотелось бы сказать, что по теореме Лебега о мажорируемой сходимости несложно подобрать c для каждой случайной величины ξ_n по отдельности, после чего взять супремум. Но он не обязательно конечен. Поэтому будем действовать так: пусть для всех $n > N$ $E[|\xi_n - \xi|] < \varepsilon$. Для всех $n \leq N$ подберём c_n такие, что $E[\xi_n[\xi_n \geq c_n]] < 2\varepsilon$, после чего возьмём максимальное из них. Поэтому конечное число членов спереди нас не пугает и достаточно рассмотреть только случай $n > N$. Заметим, что

$$E[\xi_n[\xi_n > c]] \leq E[|\xi_n - \xi|[\xi_n > c]] + E[\xi[\xi_n > c]].$$

С первым интегралом всё понятно: $E[|\xi_n - \xi|[\xi_n \geq c]] \leq E[|\xi_n - \xi|] < \varepsilon$. Для второго заметим, что

$$P(\xi_n \geq c) \leq \frac{E[\xi_n]}{c} \leq \frac{\sup_{n>N} E[\xi_n]}{c} \leq \frac{\sup_{n>N} (E[\xi] + E[|\xi_n - \xi|])}{c} \leq \frac{E[\xi] + \varepsilon}{c} \xrightarrow{c \rightarrow \infty} 0.$$

Следовательно, можно найти c такое, что можно воспользоваться абсолютной непрерывностью матожидания и получить, что $E[\xi[\xi_n \geq c]] < \varepsilon$. Тем самым получается равномерная интегрируемость.

Теперь докажем в другую сторону. Пусть последовательность ξ_n равномерно интегрируема. Введём следующую функцию:

$$f_c(x) = \begin{cases} x, & x < c \\ 0, & x \geq c+1 \\ c - c(x - c), & c \leq x < c+1 \end{cases}$$

По сути, это отсечённая функция $y = x$, дополненная линейным куском до непрерывности. Далее, заметим, что по неравенству треугольника

$$\begin{aligned} E[|\xi_n - \xi|] &= E[(\xi_n - f_c(\xi_n)) - (\xi - f_c(\xi)) + (f_c(\xi_n) - f_c(\xi))] \\ &\leq E[|\xi_n - f_c(\xi_n)|] + E[|\xi - f_c(\xi)|] + E[|f_c(\xi_n) - f_c(\xi)|] \end{aligned}$$

Теперь зафиксируем произвольный $\varepsilon > 0$ и будем ограничивать все члены по отдельности:

- Начнём с того, что при $x \geq c$ $f_c(x) \leq x$. Поэтому $E[|\xi_n - f_c(\xi_n)|] \leq E[\xi_n[\xi_n > c]]$. Но равномерная интегрируемость позволяет выбрать $c_1 = c_1(\varepsilon)$ такое, что $E[\xi_n[\xi_n > c]] < \varepsilon/3$ для всех n .
- Со вторым членом всё аналогично: выбираем $c_2 = c_2(\varepsilon)$ такое, что $E[|\xi - f_c(\xi)|] < \varepsilon/3$ для всех n .
- Для третьего достаточно вспомнить, что ξ_n сходятся к ξ по распределению и альтернативное определение сходимости по распределению. Из него следует, что $E[|f_c(\xi_n) - f_c(\xi)|] < \varepsilon/3$ для любых c и достаточно больших n .

Комбинируя эти три оценки, получим сходимость в L_1 , из которой следует сходимость матожиданий. \square

Теперь можно вернуться к нашей теореме.

Лемма. Пусть $\Phi(x)$ — функция распределения стандартной нормальной случайной величины. Тогда для всех $y \in \mathbb{R}$

$$\lim_{n \rightarrow \infty} P_{\theta_n}(T_n \geq y) = 1 - \Phi(y - \sqrt{i(\theta)}).$$

Доказательство. Для удобства будем писать, что T_n есть функция от выборки: $T_n = T_n(\mathbf{X})$. Теперь распишем $P_{\theta_n}(T_n < y)$:

$$\begin{aligned} P_{\theta_n}(T_n(\mathbf{X}) \geq y) &= \int_{\mathbf{x}: T_n(\mathbf{x}) \geq y} p_{\theta_n}(\mathbf{x}) \mu(d\mathbf{x}) = \int_{\mathbf{x}: T_n(\mathbf{x}) \geq y} \frac{p_{\theta_n}(\mathbf{x})}{p_{\theta_0}(\mathbf{x})} p_{\theta_0}(\mathbf{x}) \mu(d\mathbf{x}) \\ &= \int_{\mathbf{x}: T_n(\mathbf{x}) \geq y} \exp\{\mathcal{L}(\theta_n, \mathbf{X}) - \mathcal{L}(\theta_0, \mathbf{X})\} p_{\theta_0}(\mathbf{x}) \mu(d\mathbf{x}) \\ &= \int_{\mathbf{x}: T_n(\mathbf{x}) \geq y} \exp\left\{\sqrt{i(\theta_0)} T_n(\mathbf{x}) - \frac{i(\theta)}{2}\right\} p_{\theta_0}(\mathbf{x}) \mu(d\mathbf{x}) \\ &= e^{-i(\theta)/2} E_{\theta_0}\left[e^{T_n(\mathbf{X})\sqrt{i(\theta_0)}}[T_n(\mathbf{X}) \geq y]\right]. \end{aligned}$$

Теперь воспользуемся тем, что $T_n(\mathbf{X}) \xrightarrow{d_{\theta_0}} \xi \sim \mathcal{N}(0, 1)$ и ограниченностью функции внутри матожидания (что даёт равномерную интегрируемость). Тогда можно сказать, что

$$\mathbb{E}_{\theta_0} \left[e^{T_n(\mathbf{X})\sqrt{i(\theta_0)}} [T_n(\mathbf{X}) < y] \right] \rightarrow \mathbb{E} \left[e^{\xi\sqrt{i(\theta_0)}} [\xi < y] \right].$$

Это матожидание достаточно легко считается:

$$\begin{aligned} \mathbb{E} \left[e^{\xi\sqrt{i(\theta_0)}} [\xi < y] \right] &= \int_{-\infty}^y e^{z\sqrt{i(\theta_0)}} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz \\ &= \int_{-\infty}^y \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{z^2}{2} + z\sqrt{i(\theta_0)} - \frac{i(\theta_0)}{2} + \frac{i(\theta_0)}{2} \right\} dz \\ &= e^{i(\theta_0)/2} \int_{-\infty}^y \frac{1}{\sqrt{2\pi}} e^{-(z-\sqrt{i(\theta_0)})^2/2} dz = e^{i(\theta_0)/2} \Phi(y - \sqrt{i(\theta_0)}). \end{aligned}$$

Тем самым получаем, что

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta_n}(T_n(\mathbf{X}) < y) = \Phi(y - \sqrt{i(\theta_0)}). \quad \square$$

Теперь можно приступить к доказательству теоремы, сформулированной ещё пару страниц назад.

Доказательство. Зафиксируем $y > \sqrt{i(\theta_0)}$ и введём два события: $D_n = \{t_n > \theta_n\}$ и $S_n = \{T_n \geq y\}$. Сразу же воспользуемся леммой выше и скажем, что

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta_n}(S_n) = 1 - \Phi(y - \sqrt{i(\theta_0)}) < \frac{1}{2}.$$

Теперь заметим следующее:

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta_n}(t_n(\mathbf{X}) \leq \theta_n) = \lim_{n \rightarrow \infty} (1 - \mathbb{P}_{\theta_n}(t_n(\mathbf{X}) > \theta_n)) = 1 - \overline{\lim}_{n \rightarrow \infty} \mathbb{P}_{\theta_n}(t_n(\mathbf{X}) > \theta_n).$$

Следовательно, согласно условию

$$1 - \overline{\lim}_{n \rightarrow \infty} \mathbb{P}_{\theta_n}(t_n(\mathbf{X}) > \theta_n) \leq \frac{1}{2} \implies \overline{\lim}_{n \rightarrow \infty} \mathbb{P}_{\theta_n}(t_n(\mathbf{X}) > \theta_n) \geq \frac{1}{2}.$$

Из этого следует, что существует подпоследовательность индексов $\{n_k \mid k \in \mathbb{N}\}$ такая, что

$$\lim_{k \rightarrow \infty} \mathbb{P}_{\theta_{n_k}}(D_{n_k}) \geq \frac{1}{2} > \lim_{k \rightarrow \infty} \mathbb{P}_{\theta_{n_k}}(S_{n_k}).$$

Это означает, что при достаточно больших k $\mathbb{P}_{\theta_{n_k}}(D_{n_k}) > \mathbb{P}_{\theta_{n_k}}(S_{n_k})$. Теперь докажем ещё одну лемму:⁹

Лемма. Для всех достаточно больших k $\mathbb{P}_{\theta_0}(D_{n_k}) > \mathbb{P}_{\theta_0}(S_{n_k})$.

Доказательство. Пусть $n = n_k$ и k достаточно велико для того, чтобы выполнялось условие $\mathbb{P}_{\theta_{n_k}}(D_{n_k}) > \mathbb{P}_{\theta_{n_k}}(S_{n_k})$. Заметим, что событие S_n можно переписать следующим образом:

$$\begin{aligned} S_n = \{T_n(\mathbf{X}) \geq y\} &= \left\{ \frac{1}{\sqrt{i(\theta_0)}} \left(\mathcal{L}(\theta_n, \mathbf{X}) - \mathcal{L}(\theta_0, \mathbf{X}) + \frac{i(\theta_0)}{2} \right) > y \right\} \\ &= \{\mathcal{L}(\theta_n, \mathbf{X}) - \mathcal{L}(\theta_0, \mathbf{X}) > y'\} = \{p_{\theta_n}(\mathbf{X}) > \lambda p_{\theta_0}(\mathbf{X})\}. \end{aligned}$$

Теперь рассмотрим следующее неравенство:

$$(p_{\theta_n}(\mathbf{x}) - \lambda p_{\theta_0}(\mathbf{x}))[\mathbf{x} \in S_n] \geq (p_{\theta_n}(\mathbf{x}) - \lambda p_{\theta_0}(\mathbf{x}))[\mathbf{x} \in D_n].$$

Почему оно действительно выполняется?

- Пусть $\mathbf{x} \in S_n$. Тогда $[\mathbf{x} \in S_n] = 1$ и $[\mathbf{x} \in D_n] \leq 1$, что выполнено всегда.
- Теперь пусть $\mathbf{x} \notin S_n$. Тогда $[\mathbf{x} \in S_n] = 0$ и $p_{\theta_n}(\mathbf{x}) - \lambda p_{\theta_0}(\mathbf{x}) < 0$. Так как индикатор принимает значения 0 или 1, то получается, что неположительное число не больше нуля, что похоже на правду.

⁹Лемм много не бывает.

Проинтегрируем по мере μ :

$$P_{\theta_n}(S_n) - \lambda P_{\theta_0}(S_n) \geq P_{\theta_n}(D_n) - \lambda P_{\theta_0}(D_n)$$

Простой перегруппировкой мы получаем желаемое:

$$\lambda(P_{\theta_0}(D_n) - P_{\theta_0}(S_n)) \geq P_{\theta_n}(D_n) - P_{\theta_n}(S_n) > 0 \implies P_{\theta_0}(D_n) > P_{\theta_0}(S_n). \quad \square$$

Теперь покажем, что $P_{\theta_0}(D_n)$ и $P_{\theta_0}(S_n)$ на самом деле имеют пределы. По лемме про сходимость T_n по распределению к $\mathcal{N}(0, 1)$

$$\lim_{n \rightarrow \infty} P_{\theta_0}(S_n) = 1 - \Phi(y).$$

Далее,

$$\begin{aligned} \lim_{n \rightarrow \infty} P_{\theta_0}(D_n) &= \lim_{n \rightarrow \infty} P_{\theta_0} \left(t_n(\mathbf{X}) \geq \theta_0 + \frac{1}{\sqrt{n}} \right) = \lim_{n \rightarrow \infty} P_{\theta_0} (\sqrt{n}(t_n(\mathbf{X}) - \theta_0) \geq 1) \\ &= \lim_{n \rightarrow \infty} P_{\theta_0} \left(\frac{\sqrt{n}(t_n(\mathbf{X}) - \theta_0)}{\sigma(\theta_0)} \geq \frac{1}{\sigma(\theta_0)} \right). \end{aligned}$$

Однако по центральной предельной теореме левая часть неравенства сходится по распределению к $\mathcal{N}(0, 1)$. Следовательно, предел равен $1 - \Phi(\sigma^{-1}(\theta_0))$. Теперь воспользуемся результатом выше и скажем, что

$$\lim_{n \rightarrow \infty} P_{\theta_0}(S_n) \leq \lim_{n \rightarrow \infty} P_{\theta_0}(D_n).$$

Отсюда получаем, что

$$1 - \Phi(y) \leq 1 - \Phi\left(\frac{1}{\sigma(\theta_0)}\right) \implies \Phi\left(\frac{1}{\sigma(\theta_0)}\right) \leq \Phi(y) \implies \frac{1}{\sigma(\theta_0)} \leq y \implies \sigma^2(\theta_0) \geq \frac{1}{y^2}.$$

Однако мы брали произвольное $y > \sqrt{i(\theta_0)}$. Тогда $\sigma^2(\theta_0) \geq i^{-1}(\theta_0)$, что и требовалось доказать. \square

Теперь перейдём к самой теореме Бахадура.

Теорема 13. Пусть в условиях регулярности (R1)–(R6) для последовательности оценок $t_n(\mathbf{X})$ параметра θ выполнено условие асимптотической нормальности: для всех $\theta \in \Theta$

$$\sqrt{n}(t_n(\mathbf{X}) - \theta) \xrightarrow{d_\theta} \mathcal{N}(0, \sigma^2(\theta)).$$

Тогда для всех $\theta_0 \in \Theta$, кроме множества лебеговой меры 0, выполнено следующее условие:

$$\varliminf_{n \rightarrow \infty} P_{\theta_n}(t_n(\mathbf{X}) \leq \theta_n) \leq \frac{1}{2}.$$

Доказательство. Рассмотрим следующую функцию:

$$f_n(\theta) = \left| P_\theta(t_n(\mathbf{X}) \leq \theta) - \frac{1}{2} \right| = \left| P_\theta(\sqrt{n}(t_n(\mathbf{X}) - \theta) \leq 0) - \frac{1}{2} \right|.$$

Тогда по центральной предельной теореме для любого фиксированного $\theta \in \Theta$ $f_n(\theta) \rightarrow 0$ при $n \rightarrow \infty$. Стоит заметить, что $f_n(\theta) \in [0, 1/2]$.

Теперь введём две вещи: функцию $g_n(\theta) = f_n(\theta + n^{-1/2})$ и колмогоровскую тройку $(\mathbb{R}, \mathcal{B}(\mathbb{R}), Q)$, где Q — вероятностная мера стандартного нормального распределения. На этой тройке введём случайную величину ξ , действующую по правилу $\xi(x) = x$. Посчитаем матожидание $g_n(\xi)$:

$$E[g_n(\xi)] = \int_{-\infty}^{+\infty} g_n(x) \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \int_{-\infty}^{+\infty} f_n\left(x + \frac{1}{\sqrt{n}}\right) \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

Сделаем замену $y = x + n^{-1/2}$. Тогда

$$\begin{aligned} E[g_n(\xi)] &= \int_{-\infty}^{+\infty} f_n(y) \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(y^2 - \frac{2y}{\sqrt{n}} + \frac{1}{n}\right)\right\} dy \\ &= \int_{-\infty}^{+\infty} f_n(y) \exp\left\{\frac{y}{\sqrt{n}} - \frac{1}{2n}\right\} \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy \\ &= E\left[f_n(\xi) \exp\left\{\frac{\xi}{\sqrt{n}} - \frac{1}{2n}\right\}\right]. \end{aligned}$$

Однако

$$f_n(\xi) \exp \left\{ \frac{\xi}{\sqrt{n}} - \frac{1}{2n} \right\} \xrightarrow{\mathbf{Q}\text{-п.н.}} 0 \text{ и } f_n(\xi) \exp \left\{ \frac{\xi}{\sqrt{n}} - \frac{1}{2n} \right\} \leq \max(e^\xi, 1).$$

По теореме Лебега о мажорируемой сходимости получаем, что $E[g_n(\xi)] \rightarrow 0$. Так как $g_n(x) \in [0, 1/2]$, то это означает, что $g_n(\xi)$ сходится к нулю в L_1 . Тогда $g_n(\xi)$ сходится по вероятностной мере \mathbf{Q} к нулю. Но тогда можно выбрать подпоследовательность индексов $\{n_k, k \in \mathbb{N}\}$ такую, что $g_{n_k}(\xi) \rightarrow 0$ \mathbf{Q} -п.н.

Из этого можно сделать вывод, что существует борелевское множество $B \in \mathcal{B}(\mathbb{R})$ такое, что $\mathbf{Q}(B) = 0$ и для всех $x \notin B$ $g_{n_k}(x) \rightarrow 0$. Так как мы использовали в качестве вероятностной меры вероятностную меру стандартного нормального распределения, то это означает, что множество B имеет лебегову меру 0. Теперь заметим, что для всех $\theta \in \Theta \setminus B$

$$g_{n_k}(\theta) \rightarrow 0 \implies P_{\theta_{n_k}}(t_{n_k}(\mathbf{X}) \leq \theta_{n_k}) \rightarrow \frac{1}{2} \implies \lim_{n \rightarrow \infty} P_{\theta_n}(t_n(\mathbf{X}) \leq \theta_n) \leq \frac{1}{2}. \quad \square$$

На этом доказательство теоремы Бахадура заканчивается. Теперь можно и следствия посмотреть. Но для начала введём одно определение:

Определение 19. Пусть $\hat{\theta}_n(\mathbf{X})$ — асимптотически нормальная оценка параметра θ с асимптотической дисперсией $\sigma^2(\theta)$. Если $\sigma^2(\theta) = i^{-1}(\theta)$, то будем называть оценку $\hat{\theta}_n(\mathbf{X})$ *асимптотически эффективной*.

Следствие. В условиях регулярности оценка максимального правдоподобия асимптотически эффективна. Поэтому считается, что оценка максимального правдоподобия — это «почти» наилучшая оценка в асимптотическом плане.

Оказывается, что у оценки максимального правдоподобия есть ещё одно интересное свойство.

Теорема 14 (Эффективность оценки максимального правдоподобия). *Если в условиях неравенства Рао-Крамера $\hat{\theta}(\mathbf{X})$ является эффективной оценкой θ , то $\hat{\theta}(\mathbf{X})$ есть оценка максимального правдоподобия.*

Доказательство. Воспользуемся критерием эффективности оценки:

$$\hat{\theta}(\mathbf{X}) - \theta = \frac{1}{I_{\mathbf{X}}(\theta)} \frac{\partial}{\partial \theta} \ln p_{\theta}(\mathbf{X}) = \frac{\mathcal{L}(\theta, \mathbf{X})}{I_{\mathbf{X}}(\theta)}.$$

Так как информация Фишера положительна, то $\mathcal{L}(\theta, \mathbf{X})$ имеет тот же знак, что и $\hat{\theta}(\mathbf{X}) - \theta$: при $\theta < \hat{\theta}(\mathbf{X})$ $\mathcal{L}(\theta, \mathbf{X}) > 0$ и наоборот. Тогда получаем, что $\hat{\theta}(\mathbf{X})$ — это единственная точка максимума $f_{\theta}(\mathbf{X})$ (как функции от θ), то есть $\hat{\theta}(\mathbf{X})$ есть оценка максимального правдоподобия. \square

На этом с асимптотическим подходом заканчиваем.

4.2 Байесовский подход

Ранее мы вводили так называемый байесовский подход к сравнению оценок. В нём задача поиска наилучшей оценки параметра θ вводилась следующим образом:

$$\hat{\theta}(\mathbf{X}) = \arg \min_{\theta^*(\mathbf{X})} R_{\mathbf{Q}}(\theta^*, \theta) = \arg \min_{\theta^*(\mathbf{X})} \int_{\Theta} R(\theta^*(\mathbf{X}), t) q(t) \lambda(dt),$$

где \mathbf{Q} — вероятностная мера на Θ с плотностью $q(t)$ по мере λ . У данной плотности есть своё название:

Определение 20. Плотность $q(t)$ называют *априорной плотностью* параметра θ .

Далее, будем считать, что $\mathbf{X} = (X_1, \dots, X_n)$ — это выборка из неизвестного распределения $\mathbf{P} \in \{\mathbf{P}_{\theta} \mid \theta \in \Theta\}$, где $\{\mathbf{P}_{\theta} \mid \theta \in \Theta\}$ есть доминируемое семейство с плотностью $p_{\theta}(x)$ по мере μ . Тогда можно ввести следующую штуку:

Определение 21. *Апостериорной плотностью* параметра θ называют

$$q(t \mid \mathbf{x}) = \frac{q(t) p_t(\mathbf{x})}{\int_{\Theta} q(\tau) p_{\tau}(\mathbf{x}) \lambda(d\tau)}.$$

На данный момент будем считать, что мы работаем с квадратичной функцией потерь: $R(\hat{\theta}(\mathbf{X}), \theta) = E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2]$. Можно ли получить решение с такой функцией потерь в общем случае? Можно.

Определение 22. Байесовской оценкой параметра θ будем называть

$$\hat{\theta}_Q(\mathbf{X}) = \int_{\Theta} t q(t | \mathbf{X}) \lambda(dt).$$

На следующей лекции мы докажем, что она является наилучшей оценкой в байесовском подходе. Для этого нам понадобится следующая теорема:

Теорема 15 (о наилучшем квадратичном прогнозе). Пусть ξ и η — случайные величины. Тогда

$$\mathbb{E}[\xi | \eta] = \arg \min_{f(\eta)} \mathbb{E}[(\xi - f(\eta))^2].$$

Доказательство. Добавим и вычтем $\mathbb{E}[\xi | \eta]$, после чего раскроем:

$$\begin{aligned} \mathbb{E}[(\xi - f(\eta))^2] &= \mathbb{E}[(\xi - \mathbb{E}[\xi | \eta] + \mathbb{E}[\xi | \eta] - f(\eta))^2] \\ &= \mathbb{E}[(\xi - \mathbb{E}[\xi | \eta])^2] + 2 \mathbb{E}[(\xi - \mathbb{E}[\xi | \eta])(\mathbb{E}[\xi | \eta] - f(\eta))] + \mathbb{E}[(\mathbb{E}[\xi | \eta] - f(\eta))^2]. \end{aligned}$$

Первое слагаемое никак не зависит от f . Докажем, что второе слагаемое равно нулю. Для этого воспользуемся формулой полной вероятности и тем, что функцию от условия можно вынести за условное матожидание:

$$\begin{aligned} \mathbb{E}[(\xi - \mathbb{E}[\xi | \eta])(\mathbb{E}[\xi | \eta] - f(\eta))] &= \mathbb{E}[\mathbb{E}[(\xi - \mathbb{E}[\xi | \eta])(\mathbb{E}[\xi | \eta] - f(\eta)) | \eta]] \\ &= \mathbb{E}[(\mathbb{E}[\xi | \eta] - f(\eta)) \mathbb{E}[(\xi - \mathbb{E}[\xi | \eta]) | \eta]] \\ &= \mathbb{E}[(\mathbb{E}[\xi | \eta] - f(\eta))(\mathbb{E}[\xi | \eta] - \mathbb{E}[\xi | \eta])] = 0. \end{aligned}$$

Тем самым минимум достигается при $f(\eta) = \mathbb{E}[\xi | \eta]$. \square

5 Лекция 5

5.1 Байесовские оценки

Напомним, что в байесовском подходе задача нахождения наилучшей оценки параметра θ формулируется так:

$$\hat{\theta}(\mathbf{X}) = \arg \min_{\theta^*(\mathbf{X})} \int_{\Theta} R(\theta^*(\mathbf{X}), t) q(t) \lambda(dt),$$

где $q(t)$ — плотность вероятностного распределения Q на Θ по мере λ . Оказывается, что при некоторых условиях можно получить явное выражение для наилучшей оценки. Для доказательства этого факта сделаем несколько предположений:

- Мы будем использовать квадратичную функцию потерь: $R(\hat{\theta}(\mathbf{X}), \theta) = \mathbb{E}_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2]$.
- \mathbf{X} имеет неизвестное распределение $P \in \{P_{\theta} | \theta \in \Theta\}$, где $\{P_{\theta} | \theta \in \Theta\}$ — доминируемое семейство распределений с плотностью $p_{\theta}(x)$ по мере μ .
- Будем считать, что параметр есть число: $\Theta \in \mathbb{R}$. Далее, зафиксируем вероятностное распределение Q на Θ с плотностью $q(t)$ по мере λ .

Теперь докажем следующую теорему:

Теорема 16 (о байесовской оценке). Байесовская оценка $\hat{\theta}_Q(\mathbf{X})$ — это наилучшая оценка параметра θ в байесовском подходе.

Доказательство. Введём следующую функцию: $f(t, \mathbf{x}) = q(t) p_t(\mathbf{x})$. Заметим, что

$$\int_{\Theta \times \mathcal{X}} f(t, \mathbf{x}) \lambda(dt) \mu(d\mathbf{x}) = \int_{\Theta} q(t) \left(\int_{\mathcal{X}} p_t(\mathbf{x}) \mu(d\mathbf{x}) \right) \lambda(dt) = \int_{\Theta} q(t) \lambda(dt) = 1.$$

Тогда можно считать, что $f(t, \mathbf{x})$ есть плотность некоторого вероятностного распределения \tilde{P} на $\Theta \times \mathcal{X}$ по мере $\lambda \otimes \mu$. Далее, посмотрим на вектор (θ, \mathbf{X}) , как на случайный вектор на колмогоровской тройке $(\Theta \times \mathcal{X}, \mathcal{F}, \tilde{P})$, где \mathcal{F} — соответствующая сигма-алгебра, действующий по следующему правилу: $(\theta, \mathbf{X})(t, \mathbf{x}) = (t, \mathbf{x})$. Тогда легко заметить следующее:

$$\begin{aligned} \int_{\Theta} R(\hat{\theta}(\mathbf{X}), t) q(t) \lambda(dt) &= \int_{\Theta} \mathbb{E}_t[(\hat{\theta}(\mathbf{X}) - t)^2] q(t) \lambda(dt) \\ &= \int_{\Theta} \left(\int_{\mathcal{X}} (\hat{\theta}(\mathbf{x}) - t)^2 p_t(\mathbf{x}) \mu(d\mathbf{x}) \right) q(t) \lambda(dt) \\ &= \int_{\Theta \times \mathcal{X}} (\hat{\theta}(\mathbf{x}) - t)^2 f(t, \mathbf{x}) \lambda(dt) \mu(d\mathbf{x}). \end{aligned}$$

Но это есть ни что иное, как математическое ожидание $(\hat{\theta}(\mathbf{X}) - t)^2$ по вероятностной мере \tilde{P} . Будем обозначать его через $\tilde{E}[(\hat{\theta}(\mathbf{X}) - t)^2]$. Тогда задача поиска наилучшей оценки сводится к следующей:

$$\tilde{E}[(\hat{\theta}(\mathbf{X}) - t)^2] \rightarrow \min_{\hat{\theta}(\mathbf{X})}.$$

Теорема о наилучшем квадратичном прогнозе говорит, что ответ есть $\tilde{E}[\theta | \mathbf{X}]$. Теперь осталось показать, что $\hat{\theta}_Q(\mathbf{X}) = \tilde{E}[\theta | \mathbf{X}]$. Для этого заметим, что

- $q(t)$ есть плотность θ .
- Выборка \mathbf{X} имеет плотность

$$g(\mathbf{x}) = \int_{\Theta} f(t, \mathbf{x}) \lambda(dt).$$

- Далее, условная плотность \mathbf{X} относительно θ равна $f(t, \mathbf{x})/q(t) = p_t(\mathbf{x})$.¹⁰

Отсюда несложно получить условную плотность θ относительно \mathbf{X} . Она равна

$$\frac{f(t, \mathbf{x})}{g(\mathbf{x})} = \frac{q(t)p_t(\mathbf{x})}{\int_{\Theta} q(\tau)p_{\tau}(\mathbf{x})\lambda(d\tau)} = q(t | \mathbf{x}).$$

Тогда

$$\tilde{E}[\theta | \mathbf{X} = \mathbf{x}] = \int_{\Theta} tq(t | \mathbf{x})\lambda(dt) = \hat{\theta}_Q(\mathbf{x}) \implies \hat{\theta}_Q(\mathbf{X}) = \tilde{E}[\theta | \mathbf{X}]. \quad \square$$

Примечание. Стоит заметить, что считать условную плотность на самом деле не обязательно, так как байесовскую оценку можно представить в виде отношения интегралов по совместной плотности:

$$\hat{\theta}_Q(\mathbf{X}) = \frac{\int_{\Theta} tf(t, \mathbf{X})\lambda(dt)}{\int_{\Theta} f(t, \mathbf{X})\lambda(dt)}.$$

Теперь рассмотрим какой-нибудь пример.

Задача 7. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из равномерного распределения $U(0, \theta)$, $\theta > 0$. Найдите байесовскую оценку параметра θ , если априорная плотность распределения равна

- (a) $q(t) = [t \in [0, 1]]$,
- (b) $q(t) = t^{-2}[t \geq 1]$.

Решение. Начнём с того, что посчитаем плотность выборки:

$$p_{\theta}(\mathbf{X}) = \frac{1}{\theta^n} [0 \leq X_{(1)} \leq X_{(n)} \leq \theta] = \theta^{-n} [\theta \geq X_{(n)}].$$

- (a) Совместная плотность равна $f(t, \mathbf{X}) = q(t)p_t(\mathbf{X})$. Так как $X_{(n)} \in [0, 1]$ почти наверное, то можно написать следующее:

$$f(t, \mathbf{X}) = t^{-n} [t \in [0, 1], t \geq X_{(n)}] = t^{-n} [t \in [X_{(n)}, 1]].$$

Тогда

$$\begin{aligned} \int_{\Theta} tf(t, \mathbf{X})\lambda(dt) &= \int_{X_{(n)}}^1 t^{1-n} dt = \left. \frac{t^{2-n}}{2-n} \right|_{X_{(n)}}^1 = \frac{1}{n-2} (X_{(n)}^{2-n} - 1) \\ \int_{\Theta} f(t, \mathbf{X})\lambda(dt) &= \int_{X_{(n)}}^1 t^{-n} dt = \left. \frac{t^{1-n}}{1-n} \right|_{X_{(n)}}^1 = \frac{1}{n-1} (X_{(n)}^{1-n} - 1) \end{aligned}$$

Отсюда получаем, что байесовской оценкой в таком случае будет

$$\hat{\theta}_Q(\mathbf{X}) = \frac{n-1}{n-2} \frac{X_{(n)}^{2-n} - 1}{X_{(n)}^{1-n} - 1} = \frac{n-1}{n-2} \frac{X_{(n)} - X_{(n)}^{n-1}}{1 - X_{(n)}^{n-1}}.$$

Стоит заметить, что при больших n $\hat{\theta}_Q(\mathbf{X}) \sim X_{(n)}$.

¹⁰Из-за этого её часто обозначают $p(\mathbf{X} | \theta)$.

(b) В данном случае совместная плотность устроена немного по-другому:

$$f(t, \mathbf{X}) = t^{-(n+2)}[t \geq 1, t \geq X_{(n)}] = t^{-(n+2)}[t \geq \max(1, X_{(n)})].$$

Тогда

$$\begin{aligned} \int_{\Theta} t f(t, \mathbf{X}) \lambda(dt) &= \int_{\max(1, X_{(n)})}^{+\infty} t^{-n-1} dt = -\frac{t^{-n}}{n} \Big|_{\max(1, X_{(n)})}^{+\infty} = \frac{1}{n} (\max(1, X_{(n)}))^{-n} \\ \int_{\Theta} f(t, \mathbf{X}) \lambda(dt) &= \int_{\max(1, X_{(n)})}^{+\infty} t^{-n-2} dt = -\frac{t^{-n-1}}{n+1} \Big|_{\max(1, X_{(n)})}^{+\infty} = \frac{1}{n+1} (\max(1, X_{(n)}))^{-n-1} \end{aligned}$$

Отсюда получаем, что

$$\hat{\theta}_Q(\mathbf{X}) = \frac{n+1}{n} \max(1, X_{(n)}).$$

□

5.2 Минимаксные оценки

Помимо байесовского подхода к сравнению оценок **вводился** так называемый минимаксный подход. В нём задача поиска наилучшей оценки ставится следующим образом:

$$\hat{\theta}(\mathbf{X}) = \arg \min_{\theta^*(\mathbf{X})} \sup_{\theta \in \Theta} R(\theta^*(\mathbf{X}), \theta).$$

Сначала совершенно непонятно, можно ли вообще предъявить какой-то критерий минимаксности. Но он есть. Сделаем те же предположения, что делались для байесовских оценок:

- Мы будем использовать квадратичную функцию потерь: $R(\hat{\theta}(\mathbf{X}), \theta) = E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2]$.
- \mathbf{X} имеет неизвестное распределение $P \in \{P_{\theta} \mid \theta \in \Theta\}$, где $\{P_{\theta} \mid \theta \in \Theta\}$ — доминируемое семейство распределений с плотностью $p_{\theta}(x)$ по мере μ .
- Будем считать, что параметр есть число: $\Theta \in \mathbb{R}$. Далее, зафиксируем вероятностное распределение Q на Θ с плотностью $q(t)$ по мере λ .

Лемма (Достаточное условие минимаксности оценки). Пусть $\hat{\theta}(\mathbf{X})$ — это оценка параметра θ такое, что существует вероятностное распределение Q на Θ с плотностью $q(t)$ по мере λ , соответствующее следующему условию: для всех $\theta \in \Theta$

$$E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2] \leq \int_{\Theta} E_t[(\hat{\theta}_Q(\mathbf{X}) - t)^2] q(t) \lambda(dt).$$

Тогда $\hat{\theta}(\mathbf{X})$ есть минимаксная оценка, то есть это наилучшая оценка в минимаксном подходе.

Доказательство. Возьмём произвольную оценку $\theta^*(\mathbf{X})$. Заметим, что

$$\int_{\Theta} E_t[(\theta^*(\mathbf{X}) - t)^2] q(t) \lambda(dt) \leq \int_{\Theta} \sup_{\theta \in \Theta} E_{\theta}[(\theta^*(\mathbf{X}) - \theta)^2] q(t) \lambda(dt) = \sup_{\theta \in \Theta} E_{\theta}[(\theta^*(\mathbf{X}) - \theta)^2].$$

Однако данный интеграл можно ограничить снизу интегралом для байесовской оценки:

$$\int_{\Theta} E_t[(\hat{\theta}_Q(\mathbf{X}) - t)^2] q(t) \lambda(dt) \leq \int_{\Theta} E_t[(\theta^*(\mathbf{X}) - t)^2] q(t) \lambda(dt)$$

Тем самым получаем цепочку неравенств, которая показывает минимаксность:

$$\sup_{\theta \in \Theta} E_{\theta}[(\theta^*(\mathbf{X}) - \theta)^2] \geq \int_{\Theta} E_t[(\hat{\theta}_Q(\mathbf{X}) - t)^2] q(t) \lambda(dt) \geq \sup_{\theta \in \Theta} E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2].$$

□

Из данного утверждения можно сделать очень интересное следствие:

Следствие. Пусть для оценки $\hat{\theta}(\mathbf{X})$ существует такое вероятностное распределение Q на Θ с плотностью $q(t)$ по мере λ со следующими условиями:

- (a) Существует множество $\Psi \subseteq \Theta$ такое, что $Q(\Psi) = 1$ и для всех $\theta \in \Psi$ $E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)] = c = \text{const}$.
- (b) Для всех $\theta \notin \Psi$ $E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)] \leq c$

(с) Оценка является байесовской оценкой для распределения Q : $\hat{\theta}(X) = \hat{\theta}_Q(X)$

Тогда $\hat{\theta}(X)$ есть минимаксная оценка.

Идея доказательства состоит в том, что

$$\sup_{\theta \in \Theta} E_{\theta}[(\hat{\theta}(X) - \theta)^2] = c = \int_{\Theta} E_t[(\hat{\theta}_Q(X) - t)^2] q(t) \lambda(dt).$$

Теперь сделаем лирическое отступление. Допустим, что мы возьмём другое распределение \tilde{Q} . Что тогда можно сказать про отношение следующих интегралов:

$$\int_{\Theta} E_t[(\hat{\theta}_{\tilde{Q}}(X) - t)^2] \tilde{q}(t) \lambda(dt) \text{ и } \int_{\Theta} E_t[(\hat{\theta}_Q(X) - t)^2] q(t) \lambda(dt)?$$

На самом деле первый интеграл не больше второго. Действительно, $\hat{\theta}_{\tilde{Q}}(X)$ есть байесовская оценка для распределения \tilde{Q} . Тогда

$$\begin{aligned} \int_{\Theta} E_t[(\hat{\theta}_{\tilde{Q}}(X) - t)^2] \tilde{q}(t) \lambda(dt) &\leq \int_{\Theta} E_t[(\hat{\theta}_Q(X) - t)^2] \tilde{q}(t) \lambda(dt) \\ &\leq \sup_{\theta \in \Theta} E_{\theta}[(\hat{\theta}_Q(X) - \theta)^2] = \int_{\Theta} E_t[(\hat{\theta}_Q(X) - t)^2] q(t) \lambda(dt). \end{aligned}$$

То есть распределение Q таково, что для него значение интеграла максимально. У такого распределения есть название:

Определение 23. Вероятностное распределение Q на Θ называется *наихудшим априорным распределением*.

Теперь попробуем посчитать минимаксную оценку.

Задача 8. Пусть $X = (X_1, \dots, X_n)$ — выборка из распределения Бернулли $\text{Bern}(\theta)$. Найти минимаксную оценку для θ .

Решение. Казалось бы, совершенно непонятно, как подступаться к этой задаче, даже с учётом достаточного условия выше. Начнём с того, что попробуем найти оценку, для которой функция потерь не зависит от параметра. Для этого поэкспериментируем с обычными оценками. Начнём с великолепной оценки — с выборочного среднего \bar{X} . Для неё

$$E_{\theta}[(\bar{X} - \theta)^2] = D_{\theta}[\bar{X}] = \frac{\theta(1 - \theta)}{n}.$$

Как видно, для неё функция потерь зависит от параметра. Сделаем грязный трюк и будем искать оценку в виде $\hat{\theta}(X) = a\bar{X} + b$, где $a > 0$ и b — некоторые константы. Тогда

$$\begin{aligned} E_{\theta}[(a\bar{X} + b - \theta)^2] &= E_{\theta}[(a(\bar{X} - \theta) + (b - (1 - a)\theta))^2] \\ &= a^2 E_{\theta}[(\bar{X} - \theta)^2] + (b - (1 - a)\theta)^2 \\ &= \frac{a^2(\theta - \theta^2)}{n} + b^2 - 2b\theta(1 - a) + \theta^2(1 - a)^2 \\ &= \theta^2 \left(-\frac{a^2}{n} + (1 - a)^2 \right) + \theta \left(\frac{a^2}{n} - 2b(1 - a) \right) + b^2. \end{aligned}$$

Занулим коэффициенты перед θ и θ^2 . Тогда получается следующая система уравнений:

$$\begin{cases} n(1 - a)^2 = a^2 \\ 2nb(1 - a) = a^2 \end{cases} \implies \begin{cases} \sqrt{n} - \sqrt{na} = a \\ 2b = 1 - a \end{cases}$$

Тем самым получаем, что

$$a = \frac{\sqrt{n}}{\sqrt{n} + 1}, \quad b = \frac{1}{2(\sqrt{n} + 1)}.$$

Отсюда получаем, что в качестве кандидата можно рассматривать оценку

$$\hat{\theta}(X) = \frac{\sqrt{n}}{\sqrt{n} + 1} \bar{X} + \frac{1}{2(\sqrt{n} + 1)}.$$

Для неё функция риска будет равна

$$E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2] = \frac{1}{4(\sqrt{n} + 1)^2}.$$

Теперь попробуем найти распределение \mathbf{Q} такое, что для него $\hat{\theta}(\mathbf{X})$ будет байесовской оценкой. Для этого посмотрим на функцию правдоподобия выборки:

$$p_{\theta}(\mathbf{X}) = \prod_{i=1}^n p_{\theta}(X_i) = \prod_{i=1}^n \theta^{X_i} (1 - \theta)^{1-X_i} = \theta^{n\bar{X}} (1 - \theta)^{n(1-\bar{X})}.$$

Теперь будем смотреть на неё, как на функцию от θ . Ничего не напоминает? Правильно, это напоминает бета-распределение. Возьмём в качестве априорного распределения \mathbf{Q} бета-распределение $B(\alpha, \beta)$:

$$q(t) = \frac{t^{\alpha-1}(1-t)^{\beta-1}}{B(\alpha, \beta)}.$$

Тогда совместная плотность распределения равна

$$q(\theta)p_{\theta}(\mathbf{X}) = \frac{\theta^{n\bar{X}+\alpha-1}(1-\theta)^{n(1-\bar{X})+\beta-1}}{B(\alpha, \beta)}.$$

Так как знаменатель при подсчёте будет только зависеть от \mathbf{X} , то можно смело сказать, что $q(t|\mathbf{X})$ равно плотности бета-распределения $B(n\bar{X} + \alpha, n(1-\bar{X}) + \beta)$.¹¹ Теперь найдём математическое ожидание случайной величины $\xi \sim B(\alpha, \beta)$:

$$E[\xi] = \frac{1}{B(\alpha, \beta)} \int_0^1 t^{\alpha}(1-t)^{\beta-1} dt = \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)} = \frac{\Gamma(\alpha+1)\Gamma(\beta)\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)\Gamma(\alpha+\beta+1)} = \frac{\alpha}{\alpha+\beta}.$$

Тогда

$$\hat{\theta}_{\mathbf{Q}}(\mathbf{X}) = \frac{n\bar{X} + \alpha}{n + \alpha + \beta}.$$

Теперь мы хотим подобрать параметры α и β так, чтобы байесовская оценка $\hat{\theta}_{\mathbf{Q}}(\mathbf{X})$ была равна нашей оценке $\hat{\theta}(\mathbf{X})$. Для этого заметим, что

$$\hat{\theta}(\mathbf{X}) = \frac{2\sqrt{n}\bar{X} + 1}{2\sqrt{n} + 2} = \frac{n\bar{X} + \sqrt{n}/2}{n + n\sqrt{2}}.$$

Тогда $\alpha = \beta = \sqrt{n}/2$. Следовательно, $\mathbf{Q} = B(\sqrt{n}/2, \sqrt{n}/2)$ есть наихудшее априорное распределение и $\hat{\theta}(\mathbf{X})$ есть минимаксная оценка. \square

Иногда бывают случаи, когда наихудшее распределение найти не получается. В таких случаях спасает следующая лемма:

Лемма. Пусть $\hat{\theta}(\mathbf{X})$ — оценка параметра θ такая, что существует последовательность вероятностных мер $\{\mathbf{Q}_k, k \in \mathbb{N}\}$ на Θ с плотностями $q_k(t)$ по соответствующим мерам λ_k , удовлетворяющая следующему условию: для всех $\theta \in \Theta$

$$E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2] \leq \lim_{k \rightarrow \infty} \int_{\Theta} E_t[(\hat{\theta}_{\mathbf{Q}_k}(\mathbf{X}) - t)^2] q_k(t) \lambda_k(dt).$$

Тогда $\hat{\theta}(\mathbf{X})$ есть минимаксная оценка.

Доказательство. Доказательство почти дословно повторяет доказательство в случае одной. Зафиксируем произвольную оценку $\theta^*(\mathbf{X})$. Далее, для любого натурального k

$$\int_{\Theta} E_t[(\theta^*(\mathbf{X}) - t)^2] q_k(t) \lambda_k(dt) \leq \int_{\Theta} \sup_{\theta \in \Theta} E_{\theta}[(\theta^*(\mathbf{X}) - \theta)^2] q_k(t) \lambda_k(dt) = \sup_{\theta \in \Theta} E_{\theta}[(\theta^*(\mathbf{X}) - \theta)^2].$$

Однако данный интеграл можно ограничить снизу интегралом для байесовской оценки:

$$\int_{\Theta} E_t[(\hat{\theta}_{\mathbf{Q}_k}(\mathbf{X}) - t)^2] q_k(t) \lambda_k(dt) \leq \int_{\Theta} E_t[(\theta^*(\mathbf{X}) - t)^2] q_k(t) \lambda_k(dt)$$

Следовательно, получаем желаемое:

$$\sup_{\theta \in \Theta} E_{\theta}[(\theta^*(\mathbf{X}) - \theta)^2] \geq \lim_{k \rightarrow \infty} \int_{\Theta} E_t[(\hat{\theta}_{\mathbf{Q}_k}(\mathbf{X}) - t)^2] q_k(t) \lambda_k(dt) \geq \sup_{\theta \in \Theta} E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta)^2]. \quad \square$$

¹¹ Небольшое примечание: когда априорное $q(t)$ и апостериорное $q(t|\mathbf{X})$ распределения параметров относятся к одному типу, но могут различаться параметрами, то говорят, что распределение $q(t|\mathbf{X})$ есть сопряжённое распределению $p_t(\mathbf{X})$. В данном случае было показано, что сопряжённым для распределения Бернулли является бета-распределение.

Задача 9. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из нормального распределения $\mathcal{N}(\theta, 1)$. Найти минимаксную оценку для θ .

Решение. Для начала заметим, что оценка $\bar{\mathbf{X}}$ подходит в качестве претендента на минимаксность, так как функция потерь не зависит от параметра.

$$\mathbb{E}_\theta[(\bar{\mathbf{X}} - \theta)^2] = \mathbb{D}_\theta[\bar{\mathbf{X}}] = \frac{1}{n}.$$

Далее посмотрим на функцию правдоподобия для выборки:

$$p_\theta(\mathbf{X}) = \prod_{i=1}^n p_\theta(X_i) = \frac{1}{(2\pi)^{n/2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (X_i - \theta)^2 \right\}.$$

Если посмотреть на неё, как на функцию от θ , то выражение будет напоминать плотность нормального распределения. Скажем, что $\mathbf{Q}_k = \mathcal{N}(0, \sigma_k^2)$, где σ_k^2 есть некоторое положительное число. Тогда

$$\begin{aligned} q(t) &= \frac{1}{\sigma_k \sqrt{2\pi}} \exp \left\{ -\frac{t^2}{2\sigma_k^2} \right\}, \\ q(t)p_t(\mathbf{X}) &= \frac{1}{\sigma_k^2 (2\pi)^{(n+1)/2}} \exp \left\{ -\frac{1}{2} \left(\frac{t^2}{\sigma_k^2} + nt^2 - 2tn\bar{\mathbf{X}} + \sum_{i=1}^n X_i^2 \right) \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \left(n + \frac{1}{\sigma_k^2} \right) \left(t - \frac{n}{n + \sigma_k^{-2}} \bar{\mathbf{X}} \right)^2 \right\} \end{aligned}$$

Отсюда можно сделать вывод, что¹²

$$q(t | \mathbf{X}) \sim \mathcal{N} \left(\frac{n}{n + \sigma_k^{-2}} \bar{\mathbf{X}}, \frac{1}{n + \sigma_k^{-2}} \right).$$

Тогда

$$\hat{\theta}_{\mathbf{Q}_k}(\mathbf{X}) = \frac{n}{n + \sigma_k^{-2}} \bar{\mathbf{X}}.$$

Великолепно, мы взяли какую-то последовательность вероятностных распределений. Но можно ли сказать, что будет выполнено условие леммы? Оказывается, что можно. Заметим, что

$$\begin{aligned} \int_{\Theta} \mathbb{E}_t[(\hat{\theta}_{\mathbf{Q}_k}(\mathbf{X}) - t)^2] q_k(t) \lambda_k(dt) &= \int_{\mathbb{R}^n} \int_{\Theta} (\hat{\theta}_{\mathbf{Q}_k}(\mathbf{x}) - t)^2 q_k(t) p_t(\mathbf{x}) \lambda_k(dt) \mu(d\mathbf{x}) \\ &= \int_{\mathbb{R}^n} \int_{\Theta} (t - \hat{\theta}_{\mathbf{Q}_k}(\mathbf{x}))^2 q(t | \mathbf{x}) g(\mathbf{x}) \lambda_k(dt) \mu(d\mathbf{x}) \\ &= \int_{\mathbb{R}^n} \left(\int_{\Theta} (t - \hat{\theta}_{\mathbf{Q}_k}(\mathbf{x}))^2 q(t | \mathbf{x}) \lambda_k(dt) \right) g(\mathbf{x}) \mu(d\mathbf{x}) \\ &= \int_{\mathbb{R}^n} \frac{1}{n + \sigma_k^{-2}} g(\mathbf{x}) \mu(d\mathbf{x}) = \frac{1}{n + \sigma_k^{-2}}. \end{aligned}$$

Тогда для минимаксности $\bar{\mathbf{X}}$ должно быть выполнено следующее условие:

$$\frac{1}{n} \leq \overline{\lim}_{k \rightarrow \infty} \frac{1}{n + \sigma_k^{-2}}.$$

Но для этого достаточно взять последовательность σ_k^2 , которая будет стремиться к $+\infty$ (например, $\sigma_k^2 = k$). \square

6 Лекция 6

На этой лекции разбирались задачи из первого домашнего задания. Здесь мы оставим только формулировки.

Задача 10. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из равномерного распределения на отрезке $[0, \theta]$. Проверьте на несмещенность, состоятельность, сильную состоятельность и асимптотическую нормальность следующие оценки параметра θ : $2\bar{\mathbf{X}}$, $\bar{\mathbf{X}} + X_{(n)}/2$, $(n+1)X_{(1)}$, $X_{(1)} + X_{(n)}$, $\frac{n+1}{n} X_{(n)}$.

¹²Сопряжённое для нормального с известной дисперсией есть нормальное.

Задача 11. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из гамма-распределения с параметрами (α, λ) . Предложите асимптотически нормальную оценку $\alpha > 0$ и вычислите ее асимптотическую дисперсию, если

(a) λ известно;

(b) λ тоже неизвестно.

Задача 12. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из равномерного распределения на отрезке $[0, \theta]$. Сравните следующие оценки параметра θ в равномерном подходе с квадратичной функцией потерь:

1. $\hat{\theta}_1(\mathbf{X}) = 2\bar{X}$;

2. $\hat{\theta}_2(\mathbf{X}) = X_{(1)} + X_{(n)}$;

3. $\hat{\theta}_3(\mathbf{X}) = \frac{n+1}{n}X_{(n)}$.

Задача 13. Пусть $\theta_1^*(\mathbf{X})$ и $\theta_2^*(\mathbf{X})$ — две “почти наилучшие” оценки параметра θ в среднеквадратичном подходе (т.е. каждая из них не хуже любой другой оценки), имеющие одинаковые математические ожидания. Докажите, что тогда для любого θ они совпадают почти наверное, т.е. $\theta_1^*(\mathbf{X}) = \theta_2^*(\mathbf{X})$ P_θ -п.н.

Задача 14. $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из распределения $\mathcal{N}(\theta, 1)$, $\theta > 0$. Сравните в равномерном подходе относительно квадратичной функции потерь оценки \bar{X} и $\max(0, \bar{X})$.

Задача 15. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из гамма-распределения с плотностью

$$p_\theta(x) = \frac{2^\theta}{\Gamma(\theta)} x^{\theta-1} e^{-2x} [x \geq 0].$$

где $\theta > 0$ — неизвестный параметр. Для каких функций $\tau(\theta)$ существует эффективная оценка? Найдите информацию Фишера $i(\theta)$ одного элемента выборки.

Задача 16. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из нормального распределения с параметрами (μ, σ^2) . Найдите эффективную оценку

(a) параметра μ , если σ известно;

(b) параметра σ^2 , если μ известно.

Вычислите информацию Фишера одного наблюдения в обоих случаях. Найдите информационную матрицу в случае, когда оба параметра μ и σ^2 неизвестны.

Задача 17. Пусть \mathbf{X} — наблюдение из “регулярного” семейства $\{P_\theta \mid \theta \in \Theta\}$, $\Theta \subseteq \mathbb{R}^k$, $k > 1$. Докажите, что если $\hat{\theta}(\mathbf{X})$ — эффективная оценка θ , то она является оценкой максимального правдоподобия для θ .

Задача 18. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из распределения с плотностью

$$p_\theta(x) = \frac{1}{\alpha} \exp \left\{ -\frac{x - \beta}{\alpha} \right\} [x \geq \beta],$$

где $\theta = (\alpha, \beta)$, $\alpha > 0$ — двумерный параметр. Найдите для θ оценку максимального правдоподобия. Докажите, что полученная для α оценка $\hat{\alpha}_n(\mathbf{X})$ является асимптотически нормальной, и найдите ее асимптотическую дисперсию.

7 Лекция 7

7.1 Достаточные статистики и оптимальные оценки

Ранее мы рассмотрели асимптотический, байесовский и минимаксный подходы. Теперь осталось рассмотреть самый сильный подход — **равномерный**. Напомню, что в нём поиск наилучшей оценки параметра θ устроен так:

$$E_\theta[(\hat{\theta}(\mathbf{X}) - \theta)^2] \rightarrow \min_{\hat{\theta}(\mathbf{X})} \text{равномерно по всем } \theta \in \Theta.$$

Однако, как обсуждалось ранее, если не сузить класс оценок, то эта задача будет бессмысленной (так как риск будет должен быть равен тождественному нулю). Поэтому скажем, что мы работаем в классе *несмещённых* оценок. В таком случае $E_\theta[(\hat{\theta}(\mathbf{X}) - \theta)^2] = D_\theta[\hat{\theta}(\mathbf{X})]$. Такая переформулировка сразу приводит нас к определению оптимальной оценки.

Определение 24. Пусть $\hat{\theta}(\mathbf{X})$ — несмещённая оценка параметра $\theta \in \mathbb{R}$. Будем говорить, что $\hat{\theta}(\mathbf{X})$ есть *оптимальная* оценка, если у неё равномерно наименьшая дисперсия, то есть для любой другой несмещённой оценки $\theta^*(\mathbf{X})$

$$D_{\theta}[\hat{\theta}(\mathbf{X})] \leq D_{\theta}[\theta^*(\mathbf{X})] \text{ для всех } \theta \in \Theta.$$

Данное определение дословно переносится на многомерный случай, но там неравенство на дисперсии записывается следующим образом:

$$D_{\theta}[\hat{\theta}(\mathbf{X})] \preceq D_{\theta}[\theta^*(\mathbf{X})] \iff D_{\theta}[\theta^*(\mathbf{X})] - D_{\theta}[\hat{\theta}(\mathbf{X})] \text{ неотрицательно определена.}$$

Определение 25. Пусть $\{P_{\theta} \mid \theta \in \Theta\}$ — некоторое семейство распределений, а \mathbf{X} — наблюдение с неизвестным распределением $P \in \{P_{\theta} \mid \theta \in \Theta\}$. Будем называть статистику $S(\mathbf{X})$ *достаточной*, если существует вариант условного распределения $P_{\theta}(\mathbf{X} \in B \mid S(\mathbf{X}) = s)$, который не зависит от параметра θ , то есть существует измеримая функция $P(B, s)$ такая, что

$$P_{\theta}(\mathbf{X} \in B \mid S(\mathbf{X}) = s) = P(B, s) \text{ п.н. по распределению } S(\mathbf{X}).$$

Теорема 17 (Колмогоров-Блэкуэлл-Рао). Пусть $\hat{\theta}(\mathbf{X})$ — некоторая несмещённая оценка $\tau(\theta) \in \mathbb{R}$ с конечным матожиданием: $E_{\theta}[\hat{\theta}(\mathbf{X})] < +\infty$ для всех $\theta \in \Theta$. Далее, пусть $S(\mathbf{X})$ — это достаточная статистика для семейства распределений $\{P_{\theta} \mid \theta \in \Theta\}$. Тогда

- Пусть $\theta^*(\mathbf{X}) = E_{\theta}[\hat{\theta}(\mathbf{X}) \mid S(\mathbf{X})]$. Тогда $\theta^*(\mathbf{X})$ есть несмещённая оценка $\tau(\theta)$.
- Для всех $\theta \in \Theta$ $D_{\theta}[\theta^*(\mathbf{X})] \leq D_{\theta}[\hat{\theta}(\mathbf{X})]$.
- Равенство в неравенстве выше для всех $\theta \in \Theta$ достигается тогда и только тогда, когда $\hat{\theta}(\mathbf{X})$ является $S(\mathbf{X})$ -измеримой функцией.

Доказательство. Несмещённость $\theta^*(\mathbf{X})$ сразу же следует из формулы полной вероятности:

$$E_{\theta}[\theta^*(\mathbf{X})] = E_{\theta}[E_{\theta}[\hat{\theta}(\mathbf{X}) \mid S(\mathbf{X})]] = E_{\theta}[\hat{\theta}(\mathbf{X})] = \tau(\theta).$$

Менее очевидным фактом является то, что это действительно оценка, то есть то, что $\theta^*(\mathbf{X})$ не зависит от θ . Однако $\theta^*(\mathbf{X})$ есть интеграл $\hat{\theta}(\mathbf{X})$ по условному распределению \mathbf{X} относительно $S(\mathbf{X})$. Но ни оно, ни функция $\hat{\theta}(\mathbf{x})$ не зависят от θ . Тем самым получаем желаемое.

Теперь покажем, что ограничение на дисперсии действительно выполнено. Для этого заметим, что функция $h(x) = (x - \tau(\theta))^2$ выпукла. Тогда можно воспользоваться неравенством Йенсена:

$$(\theta^*(\mathbf{X}) - \tau(\theta))^2 = (E_{\theta}[\hat{\theta}(\mathbf{X}) \mid S(\mathbf{X})] - \tau(\theta))^2 \leq E_{\theta}[(\hat{\theta}(\mathbf{X}) - \tau(\theta))^2 \mid S(\mathbf{X})].$$

Возьмём матожидание с обеих сторон:

$$D_{\theta}[\theta^*(\mathbf{X})] \leq D_{\theta}[\hat{\theta}(\mathbf{X})].$$

Осталось понять, когда будет выполняться равенство в данном неравенстве. Для этого заметим, что

$$\begin{aligned} D_{\theta}[\hat{\theta}(\mathbf{X})] - D_{\theta}[\theta^*(\mathbf{X})] &= E_{\theta}[(\hat{\theta}(\mathbf{X}) - \tau(\theta))^2] - E_{\theta}[(\theta^*(\mathbf{X}) - \tau(\theta))^2] \\ &= E_{\theta}[(\hat{\theta}(\mathbf{X}))^2] - E_{\theta}[(\theta^*(\mathbf{X}))^2] \\ &= E_{\theta}[E_{\theta}[(\hat{\theta}(\mathbf{X}))^2 - (\theta^*(\mathbf{X}))^2 \mid S(\mathbf{X})]] \\ &= E_{\theta}[E_{\theta}[(\hat{\theta}(\mathbf{X}))^2 \mid S(\mathbf{X})] - (\theta^*(\mathbf{X}))^2]. \end{aligned}$$

Выражение внутри матожидания напоминает условную дисперсию. Действительно, по линейности условного математического ожидания

$$\begin{aligned} E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta^*(\mathbf{X}))^2 \mid S(\mathbf{X})] &= E_{\theta}[(\hat{\theta}(\mathbf{X}))^2 - 2\hat{\theta}(\mathbf{X})\theta^*(\mathbf{X}) + (\theta^*(\mathbf{X}))^2 \mid S(\mathbf{X})] \\ &= E_{\theta}[(\hat{\theta}(\mathbf{X}))^2 \mid S(\mathbf{X})] - 2\theta^*(\mathbf{X}) E_{\theta}[\hat{\theta}(\mathbf{X}) \mid S(\mathbf{X})] + (\theta^*(\mathbf{X}))^2 \\ &= E_{\theta}[(\hat{\theta}(\mathbf{X}))^2 \mid S(\mathbf{X})] - (\theta^*(\mathbf{X}))^2. \end{aligned}$$

Тогда

$$D_{\theta}[\hat{\theta}(\mathbf{X})] - D_{\theta}[\theta^*(\mathbf{X})] = E_{\theta}[E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta^*(\mathbf{X}))^2 \mid S(\mathbf{X})]] = E_{\theta}[(\hat{\theta}(\mathbf{X}) - \theta^*(\mathbf{X}))^2].$$

Но данное матожидание равно нулю тогда и только тогда, когда $\hat{\theta}(\mathbf{X}) = \theta^*(\mathbf{X})$ P_{θ} -п.н. для всех $\theta \in \Theta$. \square

Следствие. Теорема Колмогорова-Блэкуэлла-Рао верна и в многомерном случае.

Доказательство. Пусть $\tau(\theta) \in \mathbb{R}^k$. Возьмём произвольный ненулевой вектор $\mathbf{a} \in \mathbb{R}^k$. Тогда $\hat{\theta}_{\mathbf{a}}(\mathbf{X})$ есть несмещённая оценка $\langle \tau(\theta), \mathbf{a} \rangle$. Следовательно, по одномерному случаю

$$D_{\theta}[\langle \theta^*(\mathbf{X}), \mathbf{a} \rangle] \leq D_{\theta}[\langle \hat{\theta}(\mathbf{X}), \mathbf{a} \rangle].$$

Но это означает, что

$$\mathbf{a}^{\top} D_{\theta}[\theta^*(\mathbf{X})] \mathbf{a} \leq \mathbf{a}^{\top} D_{\theta}[\hat{\theta}(\mathbf{X})] \mathbf{a}.$$

Тогда $D_{\theta}[\theta^*(\mathbf{X})] \preceq D_{\theta}[\hat{\theta}(\mathbf{X})]$. Критерий равенства доказывается аналогично одномерному случаю. \square

Следствие. Пусть в условиях теоремы Колмогорова-Блэкуэлла-Рао для $\tau(\theta)$ существует единственная $S(\mathbf{X})$ -измеримая несмещённая оценка $\theta^*(\mathbf{X})$. Тогда $\theta^*(\mathbf{X})$ — оптимальная оценка $\tau(\theta)$.

Доказательство. Пусть оценка $\theta^*(\mathbf{X})$ не оптимальна. Тогда найдётся оценка $\hat{\theta}(\mathbf{X})$, которая будет лучше, то есть её дисперсия меньше. Но тогда по теореме Колмогорова-Блэкуэлла-Рао $\theta^{**}(\mathbf{X}) = E_{\theta}[\hat{\theta}(\mathbf{X}) | S(\mathbf{X})]$ будет не хуже, чем $\hat{\theta}(\mathbf{X})$. Однако и $\theta^*(\mathbf{X})$, и $\theta^{**}(\mathbf{X})$ $S(\mathbf{X})$ -измеримы. Следовательно, они равны и $\theta^*(\mathbf{X})$ не хуже, чем $\hat{\theta}(\mathbf{X})$. Противоречие с тем, что $\hat{\theta}(\mathbf{X})$ лучше, чем $\theta^*(\mathbf{X})$. Следовательно, $\theta^*(\mathbf{X})$ оптимальна. \square

Определение 26. Статистика $S(\mathbf{X})$ называется *полной* для семейства $\{P_{\theta} | \theta \in \Theta\}$, если из того, что

$$E_{\theta}[f(S(\mathbf{X}))] = 0 \text{ для любого } \theta \in \Theta$$

следует, что $f(S(\mathbf{X})) = 0$ P_{θ} -п.н. для всех $\theta \in \Theta$.

Теорема 18 (об оптимальной оценке). Пусть $S(\mathbf{X})$ — полная достаточная статистика для семейства $\{P_{\theta} | \theta \in \Theta\}$. Тогда если $\phi(S(\mathbf{X}))$ есть несмещённая оценка $\tau(\theta)$, то она будет оптимальной оценкой.

Доказательство. Пусть $\psi(S(\mathbf{X}))$ — другая несмещённая оценка $\tau(\theta)$. Тогда для любого $\theta \in \Theta$

$$E_{\theta}[\phi(S(\mathbf{X})) - \psi(S(\mathbf{X}))] = 0.$$

Но тогда по полноте $\phi(S(\mathbf{X})) = \psi(S(\mathbf{X}))$ P_{θ} -п.н. Отсюда получаем, что $\phi(S(\mathbf{X}))$ есть единственная $S(\mathbf{X})$ -измеримая несмещённая оценка $\tau(\theta)$. Следовательно, она оптимальна. \square

Тем самым, если есть полная достаточная статистика, то можно легко находить оптимальные оценки для $\tau(\theta)$, решая уравнение несмещённости

$$E_{\theta}[\phi(S(\mathbf{X}))] = \tau(\theta).$$

Для примера разберём одну задачу.

Задача 19. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из распределения Бернулли $\text{Bern}(\theta)$, $\theta \in (0, 1)$. Найти полную и достаточную статистику.

Доказательство. Ранее мы не вводили никаких методов поиска полных и достаточных статистик, поэтому придётся угадывать. Начнём с самого очевидного кандидата — выборочного среднего \bar{X} , а точнее, с суммы элементов $S(\mathbf{X}) = \sum_{i=1}^n X_i$. Проверим, что она полна. Так как распределение Бернулли дискретно, то достаточно проверить, что для любого $\mathbf{x} \in \{0, 1\}^n$ условная вероятность $P_{\theta}(\mathbf{X} = \mathbf{x} | S(\mathbf{X}) = k)$ не зависит от θ . Раскроем по определению:

$$P_{\theta}(\mathbf{X} = \mathbf{x} | S(\mathbf{X}) = k) = \frac{P_{\theta}(\mathbf{X} = \mathbf{x}, S(\mathbf{X}) = k)}{P_{\theta}(S(\mathbf{X}) = k)}$$

Теперь возникает два случая. Если $\sum_{i=1}^n x_i \neq k$, то события в вероятности в числителе несовместны и $Pr_{\theta}(\mathbf{X} = \mathbf{x} | S(\mathbf{X}) = k) = 0$. Иначе же событие $S(\mathbf{X}) = k$ включается в событие $\mathbf{X} = \mathbf{x}$ и

$$\begin{aligned} P_{\theta}(\mathbf{X} = \mathbf{x} | S(\mathbf{X}) = k) &= \frac{P_{\theta}(\mathbf{X} = \mathbf{x})}{P_{\theta}(S(\mathbf{X}) = k)} \mathbf{1}_{\sum_{i=1}^n x_i = k} \\ &= \frac{\theta^k (1 - \theta)^{n-k}}{\binom{n}{k} \theta^k (1 - \theta)^{n-k}} \mathbf{1}_{\sum_{i=1}^n x_i = k} = \binom{n}{k}^{-1} \mathbf{1}_{\sum_{i=1}^n x_i = k}. \end{aligned}$$

Тем самым получаем, что условная вероятность действительно не зависит от θ и $S(\mathbf{X})$ будет достаточной статистикой. Теперь проверим полноту. Для этого возьмём произвольную функцию f и посчитаем матожидание $E_{\theta}[f(S(\mathbf{X}))]$ по определению, пользуясь тем, что $S(\mathbf{X}) \sim \text{Bin}(n, \theta)$:

$$E_{\theta}[f(S(\mathbf{X}))] = \sum_{k=0}^n f(k) \binom{n}{k} \theta^k (1 - \theta)^{n-k}.$$

Далее допустим, что это матожидание равно нулю для любого $\theta \in (0, 1)$. Но это многочлен от θ степени n , то есть у него не может быть континуум корней, если он не есть тождественный ноль. Но есть одна проблема: отсюда не следует, что $f(k) = 0$ для всех $k \in \{1, 2, \dots, n\}$. Для этого преобразуем многочлен:

$$\sum_{k=0}^n f(k) \binom{n}{k} \theta^k (1-\theta)^{n-k} = (1-\theta)^n \sum_{k=0}^n f(k) \binom{n}{k} \left(\frac{\theta}{1-\theta}\right)^k.$$

Так как $1-\theta > 0$, то получаем, что для всех $k \in \{1, 2, \dots, n\}$ $f(k) \binom{n}{k} = 0$. Следовательно, $f(k) = 0$ для всех $k \in \{1, 2, \dots, n\}$ и $f(S(\mathbf{X})) = 0$ P_θ -п.н. для всех $\theta \in (0, 1)$. Тогда $S(\mathbf{X})$ есть полная статистика. \square

Следствие. \bar{X} есть оптимальная оценка θ в схеме Бернулли.

7.2 Критерий факторизации Неймана-Фишера

В примере выше нам сильно повезло, что мы сразу угадали полную достаточную статистику. Но как быть в общем случае? То есть возникает два относительно фундаментальных вопроса:

- Как находить достаточную статистику?
- Как проверять достаточные статистики на полноту?

На первый вопрос даёт ответ следующая теорема.

Теорема 19 (Критерий факторизации Неймана-Фишера). Пусть $\{P_\theta \mid \theta \in \Theta\}$ — это доминируемое семейство с плотностью $p_\theta(x)$. Тогда статистика $S(\mathbf{X})$ будет достаточной для данного семейства тогда и только тогда, когда существует представление $p_\theta(x)$ в виде

$$p_\theta(x) = h(x)\psi_\theta(S(x)),$$

где ψ_θ и h — некоторые неотрицательные борелевские функции.

Доказательство для дискретного случая. Для начала покажем, что из того, что $S(\mathbf{X})$ есть достаточная статистика, следует нужное представление. Для этого заметим, что

$$P_\theta(\mathbf{X} = \mathbf{x}) = P_\theta(\mathbf{X} = \mathbf{x} \mid S(\mathbf{X}) = S(\mathbf{x}))P_\theta(S(\mathbf{X}) = S(\mathbf{x})).$$

Первый множитель в произведении не зависит от θ по достаточности $S(\mathbf{X})$, а второй зависит только от θ и $S(\mathbf{X})$. Тогда

$$P_\theta(\mathbf{X} = \mathbf{x}) = h(\mathbf{x})\psi_\theta(S(\mathbf{x})),$$

где $h(\mathbf{x}) = P_\theta(\mathbf{X} = \mathbf{x} \mid S(\mathbf{X}) = S(\mathbf{x}))$, а $\psi_\theta(S(\mathbf{x})) = P_\theta(S(\mathbf{X}) = S(\mathbf{x}))$.

Теперь предположим, что существует представление $P_\theta(\mathbf{X} = \mathbf{x})$ в виде $h(\mathbf{x})\psi_\theta(S(\mathbf{x}))$. Тогда

$$\begin{aligned} P_\theta(\mathbf{X} = \mathbf{x} \mid S(\mathbf{X}) = s) &= \frac{P_\theta(\mathbf{X} = \mathbf{x}, S(\mathbf{X}) = s)}{P_\theta(S(\mathbf{X}) = s)} = \frac{P_\theta(\mathbf{X} = \mathbf{x})}{P_\theta(S(\mathbf{X}) = s)} \mathbf{1}_{S(\mathbf{x})=s} \\ &= \frac{P_\theta(\mathbf{X} = \mathbf{x})}{\sum_{\mathbf{y}: S(\mathbf{y})=S(\mathbf{x})} P_\theta(\mathbf{X} = \mathbf{y})} \mathbf{1}_{S(\mathbf{x})=s} \\ &= \frac{h(\mathbf{x})\psi_\theta(S(\mathbf{x}))}{\sum_{\mathbf{y}: S(\mathbf{y})=S(\mathbf{x})} h(\mathbf{y})\psi_\theta(S(\mathbf{y}))} \mathbf{1}_{S(\mathbf{x})=s} \\ &= \frac{h(\mathbf{x})}{\sum_{\mathbf{y}: S(\mathbf{y})=S(\mathbf{x})} h(\mathbf{y})} \mathbf{1}_{S(\mathbf{x})=s} \end{aligned}$$

Тем самым получаем, что условная вероятность не зависит от θ . Следовательно, $S(\mathbf{X})$ есть достаточная статистика. \square

Заметьте, что это было доказательство для дискретного случая. Доказательство для непрерывного случая не такое тривиальное.¹³ Для него требуется теорема о пересчёте условных математических ожиданий. Но для того, чтобы её сформулировать, нужно ввести несколько вещей.

Определение 27. Пусть (Ω, \mathcal{F}, P) — вероятностное пространство. Далее, пусть Q — это другая вероятностная мера на измеримом пространстве (Ω, \mathcal{F}) . Будем говорить, что вероятностная мера Q абсолютно непрерывна относительно вероятностной меры P , если для любого $A \in \mathcal{F}$ такого, что $P(A) = 0$, выполнено, что $Q(A) = 0$. Обозначение: $Q \ll P$.

¹³Владимир Васильевич, ну не аналогично они доказываются, не аналогично!

Пример. Пусть $(\Omega, \mathcal{F}) = (\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Далее, пусть $P = \mathcal{N}(0, 1)$, а $Q = \text{Exp}(1)$. Несложно понять, что $Q \ll P$. Однако в другую сторону это неверно, так как если взять $A = (-\infty, 0]$, то $Q(A) = 0$, но $P(A) = 1/2 \neq 0$.

Теорема 20 (Радон, Никодим). Пусть (Ω, \mathcal{F}) есть измеримое пространство, а P и \tilde{P} — две вероятностные меры на нём, причём $\tilde{P} \ll P$. Тогда существует единственная P -н.п. случайная величина ξ такая, что для любого $A \in \mathcal{F}$

$$\tilde{P}(A) = E[\xi \mathbf{1}_{\xi \in A}],$$

где E означает взятие математического ожидания по вероятностной мере P . Такую случайную величину называют производной Радона-Никодима и обозначают

$$\xi = \frac{d\tilde{P}}{dP}.$$

Если кому интересно доказательство теоремы, то обращайтесь к учебнику по функциональному анализу (к тому же Колмогорову-Фомину).

Теперь докажем один достаточно простой факт.

Теорема 21. Пусть ξ — это произвольная случайная величина, а $\tilde{P} \ll P$ — две вероятностные меры на измеримом пространстве (Ω, \mathcal{F}) . Тогда

$$\tilde{E}[\xi] = E \left[\xi \frac{d\tilde{P}}{dP} \right],$$

где \tilde{E} означает, что мы берём математическое ожидание по вероятностной мере \tilde{P} .

Доказательство. Пусть $\xi(\omega) = \mathbf{1}_{\omega \in A}$, где $A \in \mathcal{F}$. Тогда по теореме Радона-Никодима

$$\tilde{E}[\mathbf{1}_A] = \tilde{P}(A) = E \left[\mathbf{1}_A \frac{d\tilde{P}}{dP} \right].$$

То есть условие теоремы выполнено для индикаторов. Отсюда по линейности математического ожидания несложно получить это утверждение для простых случайных величин. Далее возьмём последовательность простых случайных величин ξ_n , монотонно приближающих ξ . Пользуясь теоремой о монотонной сходимости, получим желаемое. \square

Теорема 22 (Пересчёт условных математических ожиданий). Пусть $\tilde{P} \ll P$ — две вероятностные меры на измеримом пространстве (Ω, \mathcal{F}) . Далее, пусть ξ и η — случайные величины, причём $E[|\xi|] < +\infty$. Тогда

$$\tilde{E}[\xi | \eta] = \frac{E \left[\xi \frac{d\tilde{P}}{dP} \middle| \eta \right]}{E \left[\frac{d\tilde{P}}{dP} \middle| \eta \right]}$$

Доказательство. Покажем, что правая часть равенства удовлетворяет свойствам условного математического ожидания. Действительно, она есть отношение измеримых функций от η , то есть это измеримая функция от η . Далее проверим интегральное свойство: для произвольного борелевского множества B

$$\begin{aligned} \tilde{E} \left[\frac{E \left[\xi \frac{d\tilde{P}}{dP} \middle| \eta \right]}{E \left[\frac{d\tilde{P}}{dP} \middle| \eta \right]} \mathbf{1}_{\eta \in B} \right] &= E \left[\frac{E \left[\xi \frac{d\tilde{P}}{dP} \middle| \eta \right]}{E \left[\frac{d\tilde{P}}{dP} \middle| \eta \right]} \mathbf{1}_{\eta \in B} \frac{d\tilde{P}}{dP} \right] \\ &= E \left[E \left[\frac{E \left[\xi \frac{d\tilde{P}}{dP} \middle| \eta \right]}{E \left[\frac{d\tilde{P}}{dP} \middle| \eta \right]} \mathbf{1}_{\eta \in B} \frac{d\tilde{P}}{dP} \middle| \eta \right] \right] \\ &= E \left[E \left[\xi \frac{d\tilde{P}}{dP} \middle| \eta \right] \mathbf{1}_{\eta \in B} \right] = E \left[\xi \frac{d\tilde{P}}{dP} \mathbf{1}_{\eta \in B} \right] = \tilde{E}[\xi \mathbf{1}_{\eta \in B}]. \end{aligned}$$

Тем самым интегральное свойство выполнено и правая часть равенства действительно есть условное математическое ожидание. \square

В доказательстве был опущен один момент: делить на ноль нельзя, а проверка на то, что знаменатель не ноль, была пропущена. Но несложно показать, что он действительно не ноль.

Задача 20. Докажите, что

$$\tilde{P} \left(E \left[\frac{d\tilde{P}}{dP} \middle| \eta \right] = 0 \right) = 0.$$

Доказательство. Для удобства скажем, что $\xi = \frac{d\tilde{P}}{dP}$. Далее, по формуле полной вероятности:

$$\tilde{P}(E[\xi | \eta] = 0) = \tilde{E}[\mathbf{1}_{E[\xi | \eta]=0}] = E[\xi \mathbf{1}_{E[\xi | \eta]=0}] = E[E[\xi \mathbf{1}_{E[\xi | \eta]=0} | \eta]] = E[E[\xi | \eta] \mathbf{1}_{E[\xi | \eta]=0}] = 0. \quad \square$$

8 Лекция 8

8.1 Критерий факторизации Неймана-Фишера: продолжение

На предыдущей лекции был сформулирован критерий достаточности статистики $S(\mathbf{X})$, называемый критерием факторизации Неймана-Фишера. Но он был доказан только в дискретном случае, так как для доказательства в непрерывном случае требовалась формула пересчёта условных матожиданий. Теперь она есть и доказана, поэтому можно приступить и к критерию.

Доказательство, данное ниже, будет корректно для всех случаев, но мы сконцентрируемся на непрерывном случае.

Доказательство. Зафиксируем некоторое $\theta' \in \Theta$ и для любого $\theta \neq \theta'$ введём следующую вероятностную меру:

$$P_{\theta, \theta'} = \frac{P_\theta + P_{\theta'}}{2}.$$

Сразу же заметим, что $P_\theta \ll P_{\theta, \theta'}$ и $P_{\theta'} \ll P_{\theta, \theta'}$. Тогда по теореме Радона-Никодима существуют производные

$$\begin{aligned} \frac{dP_\theta}{dP_{\theta, \theta'}}(x) &= \frac{2p_\theta(x)}{p_\theta(x) + p_{\theta'}(x)} \equiv f_\theta(x), \\ \frac{dP_{\theta'}}{dP_{\theta, \theta'}}(x) &= \frac{2p_{\theta'}(x)}{p_\theta(x) + p_{\theta'}(x)} \equiv f_{\theta'}(x). \end{aligned}$$

Теперь выразим плотность $p_\theta(x)$ через $f_\theta(x)$:

$$p_\theta(x) = p_{\theta'}(x) \frac{f_\theta(x)}{2 - f_\theta(x)}.$$

Если доказать, что $f_\theta(x)$ есть функция от $S(x)$ то мы сразу же получим, что $p_\theta(x) = h(x)\psi_\theta(S(x))$, где $h(x) = p_{\theta'}(x)$, а $\psi_\theta(S(x)) = f_\theta(x)/(2 - f_\theta(x))$. Это план доказательства в прямую сторону.

Какой же может быть план доказательства в обратную сторону? Пусть имеет место факторизация: $p_\theta(x) = h(x)\psi_\theta(S(x))$. Тогда

$$f_\theta(x) = \frac{2h(x)\psi_\theta(S(x))}{h(x)\psi_\theta(S(x)) + h(x)\psi_{\theta'}(S(x))} = \frac{2\psi_\theta(S(x))}{\psi_\theta(S(x)) + \psi_{\theta'}(S(x))}.$$

Следовательно, $f_\theta(x)$ есть функция от $S(x)$ и θ . После чего этим можно будет воспользоваться.

Для начала докажем в обратную сторону. Пусть имеет место факторизация. Тогда $f_\theta(x)$ и $f_{\theta'}(x)$ есть функции от $S(x)$. Воспользуемся формулой пересчёта условного математического ожидания $T(\mathbf{X})$ по $S(\mathbf{X})$:

$$\begin{aligned} E_\theta[T(\mathbf{X}) | S(\mathbf{X})] &= \frac{E_{\theta, \theta'}[T(\mathbf{X})f_\theta(\mathbf{X}) | S(\mathbf{X})]}{E_{\theta, \theta'}[f_\theta(\mathbf{X}) | S(\mathbf{X})]} = \frac{f_\theta(\mathbf{X}) E_{\theta, \theta'}[T(\mathbf{X}) | S(\mathbf{X})]}{f_\theta(\mathbf{X})} \\ &= E_{\theta, \theta'}[T(\mathbf{X}) | S(\mathbf{X})]. \end{aligned}$$

Аналогично получаем, что $E_{\theta'}[T(\mathbf{X}) | S(\mathbf{X})] = E_{\theta, \theta'}[T(\mathbf{X}) | S(\mathbf{X})]$. Следовательно, для любого $\theta \neq \theta'$ $E_\theta[T(\mathbf{X}) | S(\mathbf{X})] = E_{\theta'}[T(\mathbf{X}) | S(\mathbf{X})]$. Теперь возьмём $T(\mathbf{X}) = \mathbf{1}_{\mathbf{X} \in B}$, где $B \in \mathcal{B}(\mathbb{R}^n)$. Отсюда получаем, что условное распределение \mathbf{X} относительно $S(\mathbf{X})$ не зависит от θ , но это и означает достаточность статистики $S(\mathbf{X})$.

Приступим к доказательству в прямую сторону. Нам нужно доказать, что $f_\theta(x)$ есть $S(x)$ -измеримая функция. Для этого покажем, что $E_\theta[f_\theta(x) | S(\mathbf{X})] = f_\theta(\mathbf{X})$. В силу достаточности статистики $S(\mathbf{X})$ $E_\theta[f_\theta(x) | S(\mathbf{X})]$ не зависит от θ и

$$E_{\theta'}[f_\theta(x) | S(\mathbf{X})] = E_\theta[f_\theta(x) | S(\mathbf{X})].$$

Покажем, что $E_{\theta, \theta'}[f_\theta(x) | S(\mathbf{X})] = E_\theta[f_\theta(x) | S(\mathbf{X})]$. Для этого проверим интегральное свойство. Зафиксируем произвольное борелевское множество $B \in \mathcal{B}(\mathbb{R}^k)$ (будем считать, что $S(\mathbf{X}) \in \mathbb{R}^k$) и распишем:

$$E_{\theta, \theta'}[f_\theta(x) \mathbf{1}_{S(\mathbf{X}) \in B}] = \frac{1}{2} E_\theta[f_\theta(x) \mathbf{1}_{S(\mathbf{X}) \in B}] + \frac{1}{2} E_{\theta'}[f_\theta(x) \mathbf{1}_{S(\mathbf{X}) \in B}].$$

Далее, по интегральному свойству

$$\begin{aligned}\mathbb{E}_{\theta, \theta'}[f_{\theta}(x) \mathbf{1}_{S(\mathbf{X}) \in B}] &= \frac{1}{2} \mathbb{E}_{\theta}[\mathbb{E}_{\theta}[f_{\theta}(x) | S(\mathbf{X})] \mathbf{1}_{S(\mathbf{X}) \in B}] + \frac{1}{2} \mathbb{E}_{\theta'}[\mathbb{E}_{\theta'}[f_{\theta}(x) | S(\mathbf{X})] \mathbf{1}_{S(\mathbf{X}) \in B}] \\ &= \frac{1}{2} \mathbb{E}_{\theta}[\mathbb{E}_{\theta}[f_{\theta}(x) | S(\mathbf{X})] \mathbf{1}_{S(\mathbf{X}) \in B}] + \frac{1}{2} \mathbb{E}_{\theta'}[\mathbb{E}_{\theta'}[f_{\theta}(x) | S(\mathbf{X})] \mathbf{1}_{S(\mathbf{X}) \in B}] \\ &= \mathbb{E}_{\theta, \theta'}[\mathbb{E}_{\theta}[f_{\theta}(x) | S(\mathbf{X})] \mathbf{1}_{S(\mathbf{X}) \in B}].\end{aligned}$$

Тем самым получаем желаемое. Но вопрос: а зачем нам это нужно? Затем, чтобы применить формулу пересчёта условных математических ожиданий. Применим её к статистике $f_{\theta}(\mathbf{X})$:

$$\mathbb{E}_{\theta}[f_{\theta}(x) | S(\mathbf{X})] = \frac{\mathbb{E}_{\theta, \theta'}[f_{\theta}^2(x) | S(\mathbf{X})]}{\mathbb{E}_{\theta, \theta'}[f_{\theta}(x) | S(\mathbf{X})]}$$

Но тогда мы получаем интересную вещь:

$$\mathbb{E}_{\theta, \theta'}[f_{\theta}^2(x) | S(\mathbf{X})] = (\mathbb{E}_{\theta, \theta'}[f_{\theta}(x) | S(\mathbf{X})])^2.$$

Однако нам известно, что для любых случайных величин ξ и η $D[\xi | \eta] = \mathbb{E}[(\xi - \mathbb{E}[\xi | \eta])^2 | \eta] = \mathbb{E}[\xi^2 | \eta] - (\mathbb{E}[\xi | \eta])^2$. Тогда

$$\mathbb{E}_{\theta, \theta'}[(f_{\theta}(x) - \mathbb{E}_{\theta, \theta'}[f_{\theta}(x) | S(\mathbf{X})])^2 | S(\mathbf{X})] = 0.$$

Взяв матожидание $\mathbb{E}_{\theta, \theta'}$, получим, что $\mathbb{E}_{\theta, \theta'}[f_{\theta}(x) | S(\mathbf{X})] = f_{\theta}(x)$. Следовательно, $\mathbb{E}_{\theta}[f_{\theta}(x) | S(\mathbf{X})] = f_{\theta}(\mathbf{X})$ и получается искомое разложение $p_{\theta}(x)$, что и требовалось доказать. \square

Попробуем применить этот критерий. Ранее мы находили оптимальную оценку для схемы Бернулли. Теперь найдём её для равномерного распределения.

Задача 21. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из равномерного распределения $U[0, \theta]$, $\theta > 0$. Найти оптимальную оценку θ .

Доказательство. План действий следующий: нужно найти достаточную статистику, проверить её на полноту, после чего решить уравнение несмещённости. Первое сделать совсем несложно, пользуясь критерием факторизации. Распишем плотность выборки:

$$p_{\theta}(\mathbf{X}) = \prod_{k=1}^n p_{\theta}(X_k) = \prod_{k=1}^n \frac{1}{\theta} \mathbf{1}_{0 \leq X_k \leq \theta} = \frac{1}{\theta^n} \mathbf{1}_{0 \leq X_{(1)} \leq X_{(n)} \leq \theta}.$$

Достаточной статистикой в данном случае будет являться $X_{(n)}$, так как в критерии факторизации можно взять $h(\mathbf{X}) = \mathbf{1}_{0 \leq X_{(1)} \leq X_{(n)}}$, $\psi_{\theta}(S(\mathbf{X})) = \theta^{-n} \mathbf{1}_{X_{(n)} \leq \theta}$.

Небольшое лирическое отступление: а почему мы рассматриваем именно эту статистику? Достаточных статистик полно — как минимум, сама выборка будет достаточной. Но есть такая «рекомендация»: если параметр k -мерный, то статистику тоже подбирают k -мерной. Но так можно сделать не всегда. Например, если взять $\mathcal{N}(\mu, \sigma^2)$, то достаточной статистикой будет $(\sum_{k=1}^n X_k, \sum_{k=1}^n X_k^2)$. Однако, если взять распределение $\mathcal{N}(\theta^2, \theta^2)$, то связь между μ и σ^2 есть, а связь между $\sum_{k=1}^n X_k$ и $\sum_{k=1}^n X_k^2$ построить особо не получится. Тем самым получается, что для одномерного параметра θ статистика $(\sum_{k=1}^n X_k, \sum_{k=1}^n X_k^2)$ будет двумерной достаточной, но она не будет полной.

Вернёмся к задаче и покажем, что $X_{(n)}$ есть полная статистика. Для этого вспомним, что плотность максимума равна

$$p_{X_{(n)}}(x) = \frac{nx^{n-1}}{\theta^n} \mathbf{1}_{x \in [0, \theta]}.$$

Посчитаем матожидание произвольной функции $g(X_{(1)})$:

$$\mathbb{E}_{\theta}[g(X_{(n)})] = \int_0^{\theta} g(x) \frac{nx^{n-1}}{\theta^n} dx.$$

Теперь предположим, что для всех $\theta > 0$ $\mathbb{E}_{\theta}[g(X_{(n)})] = 0$. Это равносильно тому, что

$$\int_0^{\theta} g(x) x^{n-1} dx = 0.$$

Продифференцируем по θ :

$$g(\theta) \theta^{n-1} = 0 \implies g(\theta) = 0 \text{ для всех } \theta > 0.$$

Тогда $g(X_{(1)}) = 0$ P_θ -п.н. для всех $\theta > 0$ и $X_{(n)}$ есть полная достаточная статистика. Осталось решить уравнение несмещённости: $E_\theta[\phi(X_{(1)})] = \theta$. Заметим, что это равносильно тому, что:

$$\frac{\theta^{n+1}}{n} = \int_0^\theta \phi(x) x^{n-1} dx.$$

Продифференцируем по θ :

$$\frac{n+1}{n} \theta^n = \phi(\theta) \theta^{n-1} \implies \phi(\theta) = \frac{n+1}{n} \theta.$$

Следовательно, оптимальной оценкой θ является $(1 + n^{-1})X_{(n)}$. \square

Что ещё бывает здесь полезного? Например, можно считать условные математические ожидания через полные достаточные статистики. Рассмотрим какой-нибудь плохой пример.

Задача 22. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из стандартного равномерного распределения $U[0, 1]$. Вычислить $E[\sqrt{X_1 + X_2} | X_{(n)}]$.

Доказательство. Как бы мы могли считать такую вещь? Можно было бы попытаться посчитать условную плотность, но никто не гарантирует её существования, да и функции очень плохие — вычисления будут отвратительными. Попробуем свести её к статистической задаче. Для этого скажем, что мы смотрим не на стандартное равномерное распределение, а на $U[0, \theta]$, где $\theta > 0$ (потом просто подставим $\theta = 1$). В таком случае известно, что $X_{(n)}$ есть полная достаточная статистика в равномерной модели $U[0, \theta]$. Тогда $E_\theta[\sqrt{X_1 + X_2} | X_{(n)}]$ не зависит от θ , а зависит только от $X_{(n)}$. И чем же это является? Ответ такой: оптимальной оценкой своего математического ожидания, то есть

$$E_\theta[E_\theta[\sqrt{X_1 + X_2} | X_{(n)}]] = E_\theta[\sqrt{X_1 + X_2}].$$

Осталось посчитать это матожидание:

$$\begin{aligned} E_\theta[\sqrt{X_1 + X_2}] &= \int_0^\theta \int_0^\theta \sqrt{x+y} \frac{1}{\theta^2} dx dy = \frac{1}{\theta^2} \int_0^\theta \left(\int_0^\theta \sqrt{x+y} dx \right) dy \\ &= \frac{1}{\theta^2} \int_0^\theta \frac{2}{3} (y + \theta)^{3/2} dy = \frac{1}{\theta^2} \frac{2}{3} \frac{2}{5} (2\theta)^{5/2} = \frac{16\sqrt{2}}{15} \sqrt{\theta}. \end{aligned}$$

Тогда получаем, что

$$E_\theta[\sqrt{X_1 + X_2} | X_{(n)}] = \frac{16\sqrt{2}}{15} \sqrt{X_{(n)}}. \quad \square$$

Этот метод позволяет быстро считать условные матожидания, но он работает только тогда, когда в условии стоит полная достаточная статистика. На этом с достаточностью заканчиваем.

8.2 Полнота

Если мы ищем достаточную статистику, то мы параллельно проверяем её на полноту. Но как это делать? В случае схемы Бернулли или же равномерного проверки были сделана честно, но в общем случае это достаточно трудно.

Это достаточно тяжёлый вопрос. Для достаточности есть достаточно общий критерий факторизации, который даёт ответ в большинстве разумных случаев. Для полноты же такого удобного критерия нет. Однако в некоторых хороших случаях можно вывести условия для полноты — например, если семейство распределений является экспоненциальным.

Определение 28. Пусть $\{P_\theta \mid \theta \in \Theta\}$ — семейство распределений с k -мерным параметром: $\Theta \subseteq \mathbb{R}^k$, $k \in \mathbb{N}$. Далее, пусть оно доминируемо с плотностью $p_\theta(x)$ по мере μ . Если эта плотность представима в следующем виде:

$$p_\theta(x) = h(x) \exp \left\{ \sum_{i=1}^k u_i(x) a_i(\theta) + v(\theta) \right\},$$

где h и u_i — это борелевские функции, то семейство распределений $\{P_\theta \mid \theta \in \Theta\}$ называется *экспоненциальным*.

Ранее мы уже видели примеры семейств распределений, не являющихся экспоненциальным — то же равномерное, так как индикатор в виде экспоненты не представить.

Теперь предположим, что семейство является экспоненциальным. Чему будет равна достаточная статистика? Понятно, что она будет равна $\mathbf{S}(x) = (u_1(x), \dots, u_k(x))$. Оказывается, что если функции $a_i(\theta)$ достаточно хороши, то $\mathbf{u}(x)$ будет и полной.

Теорема 23 (об экспоненциальном семействе). Пусть $\{P_\theta \mid \theta \in \Theta\}$ — экспоненциальное семейство распределений, а функция $\mathbf{a}(\theta) = (a_1(\theta), \dots, a_k(\theta))$ такова, что при пробегании θ всего Θ $\mathbf{a}(\theta)$ зачерчивает k -мерный параллелепипед (или шар). Тогда $\mathbf{S}(\mathbf{X}) = (u_1(\mathbf{X}), \dots, u_k(\mathbf{X}))$ будет полной достаточной статистикой для θ .

Она обычно формулируется в таком виде, но есть вопрос: а важно ли то, что размерности параметра и статистики совпадают? В принципе, необязательно, но это достаточно разумное требование, так как иначе могут возникать не самые приятные подводные камни.

Доказательство данной теоремы требует две леммы, в доказательстве которых используется комплексный анализ. Первая из них описывает плотность $S(\mathbf{X})$, а вторая обычно называется теоремой единственности аналитической функции комплексного переменного. Начнём с плотности. Было бы неплохо сказать что-то про плотность $S(\mathbf{X})$, так как при рассуждении про полноту требуется распределение этой статистики.

Пусть

$$\psi_\theta(S(\mathbf{X})) = \exp \left\{ \sum_{i=1}^k u_i(x) a_i(\theta) + v(\theta) \right\},$$

где $\mathbf{S}(\mathbf{X}) = (u_1(\mathbf{X}), \dots, u_k(\mathbf{X}))$, то есть $p_\theta(\mathbf{X}) = h(\mathbf{X})\psi_\theta(\mathbf{S}(\mathbf{X}))$.

Лемма. Статистика $\mathbf{S}(\mathbf{X})$ имеет плотность $\psi_\theta(\mathbf{s})$ по мере

$$\nu(B) = \int_{S^{-1}(B)} h(\mathbf{x}) d\mathbf{x} \text{ (если } p_\theta(\mathbf{x}) \text{ абсолютно непрерывна)}.$$

Доказательство. Для существования плотности нужно показать, что распределение статистики $\mathbf{S}(\mathbf{X})$ абсолютно непрерывно относительно меры ν . Пусть G_θ есть распределение $\mathbf{S}(\mathbf{X})$. Если $\nu(B) = 0$, то $h(\mathbf{x}) = 0$ для всех $\mathbf{x} \in S^{-1}(B)$ (так как $h(\mathbf{x}) \geq 0$). Тогда посмотрим на $P_\theta(\mathbf{S}(\mathbf{X}) \in B)$. Заметим, что

$$P_\theta(\mathbf{S}(\mathbf{X}) \in B) = E_\theta[\mathbf{1}_{\mathbf{S}(\mathbf{X}) \in B}] = \int_{\mathbb{R}^k} \mathbf{1}_{\mathbf{S}(\mathbf{x}) \in B} p_\theta(\mathbf{x}) d\mathbf{x} = \int_{S^{-1}(B)} h(\mathbf{x}) \psi_\theta(\mathbf{S}(\mathbf{x})) d\mathbf{x} = 0.$$

Следовательно, распределение G_θ абсолютно непрерывно относительно ν и по теореме Радона-Никодима существует производная $\frac{dG_\theta}{d\nu}$. Осталось проверить, что она равна $\psi_\theta(\mathbf{s})$. Снова распишем $P_\theta(\mathbf{S}(\mathbf{X}) \in B)$:

$$P_\theta(\mathbf{S}(\mathbf{X}) \in B) = \int_{S^{-1}(B)} h(\mathbf{x}) \psi_\theta(\mathbf{S}(\mathbf{x})) d\mathbf{x}.$$

Теперь пересчитываем интеграл, делая замену $\mathbf{s} = \mathbf{S}(\mathbf{x})$:

$$P_\theta(\mathbf{S}(\mathbf{X}) \in B) = \int_B \psi_\theta(\mathbf{s}) \nu(d\mathbf{s}).$$

□

Вторая лемма утверждает следующее:

Лемма (единственности аналитических функций). Пусть G_1 и G_2 — две меры Лебега такие, что существует параллелепипед $D \subseteq \mathbb{R}^k$ со следующим условием: для любого $\mathbf{a} \in D$

$$\int_{\mathbb{R}^k} e^{\langle \mathbf{a}, \mathbf{y} \rangle} G_1(d\mathbf{y}) = \int_{\mathbb{R}^k} e^{\langle \mathbf{a}, \mathbf{y} \rangle} G_2(d\mathbf{y}) < \infty$$

Тогда $G_1 = G_2$.

Идея доказательства. <здесь была идея, но её я допишу позже>

□

Теперь докажем теорему об экспоненциальном семействе.

Доказательство. Достаточность статистики $\mathbf{S}(\mathbf{X})$ следует из критерия факторизации Неймана-Фишера. Покажем, что она будет полна. Пусть $E_\theta[f(\mathbf{S}(\mathbf{X}))] = 0$ для всех $\theta \in \Theta$. Заметим, что по лемме о плотности статистики

$$E_\theta[f(\mathbf{S}(\mathbf{X}))] = \int_{\mathbb{R}^k} f(\mathbf{s}) \psi_\theta(\mathbf{s}) \nu(d\mathbf{s}) = 0.$$

Теперь введём две функции: $f^+(\mathbf{s}) = \max\{f(\mathbf{s}), 0\}$ и $f^-(\mathbf{s}) = \min\{-f(\mathbf{s}), 0\}$. Тогда для любого $\theta \in \Theta$

$$\int_{\mathbb{R}^k} f^+(\mathbf{s}) \psi_\theta(\mathbf{s}) \nu(d\mathbf{s}) = \int_{\mathbb{R}^k} f^-(\mathbf{s}) \psi_\theta(\mathbf{s}) \nu(d\mathbf{s}).$$

Теперь вспомним, чему равно $\psi_\theta(\mathbf{s})$. Тогда

$$\int_{\mathbb{R}^k} f^+(\mathbf{s}) \exp \{ \langle \mathbf{a}(\boldsymbol{\theta}), \mathbf{s} \rangle + v(\boldsymbol{\theta}) \} \nu(d\mathbf{s}) = \int_{\mathbb{R}^k} f^-(\mathbf{s}) \exp \{ \langle \mathbf{a}(\boldsymbol{\theta}), \mathbf{s} \rangle + v(\boldsymbol{\theta}) \} \nu(d\mathbf{s})$$

Следовательно, существует некоторый параллелепипед $D \subseteq \mathbb{R}^k$ такой, что для любого $\mathbf{a} \in D$:

$$\int_{\mathbb{R}^k} e^{\langle \mathbf{a}, \mathbf{s} \rangle} f^+(\mathbf{s}) \nu(d\mathbf{s}) = \int_{\mathbb{R}^k} e^{\langle \mathbf{a}, \mathbf{s} \rangle} f^-(\mathbf{s}) \nu(d\mathbf{s}).$$

Но это есть ни что иное, как лемма о единственности аналитической функции. Тогда

$$f^+(\mathbf{s}) \nu(d\mathbf{s}) = f^-(\mathbf{s}) \nu(d\mathbf{s}) \implies f^+(\mathbf{s}) = f^-(\mathbf{s}) \text{ по мере } \nu.$$

Так как $\mathbf{S}(\mathbf{X})$ имеет плотность по мере ν , то $f(\mathbf{S}(\mathbf{X})) = 0$ \mathbf{P}_θ -п.н. для любого $\boldsymbol{\theta} \in \Theta$. Следовательно, $\mathbf{S}(\mathbf{X})$ полна. \square

Теперь рассмотрим пример применения теоремы.

Задача 23. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из нормального распределения $\mathcal{N}(\mu, \sigma^2)$. Рассмотрим три случая:

- (a) Среднее μ неизвестно, дисперсия σ^2 известна: $\mu = \theta$, $\theta \in \mathbb{R}$;
- (b) Среднее μ известно, дисперсия σ^2 неизвестна: $\sigma^2 = \theta$, $\theta \in \mathbb{R}_{++}$;
- (c) Ни среднее μ , ни дисперсия σ^2 неизвестны: $\boldsymbol{\theta} = (\mu, \sigma^2)$, $\boldsymbol{\theta} \in \mathbb{R} \times \mathbb{R}_{++}$.

Найти оптимальные оценки параметра θ .

Доказательство. Для начала выпишем плотность выборки:

$$p_\theta(\mathbf{X}) = \prod_{k=1}^n p_\theta(X_k) = \prod_{k=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{(X_k - \mu)^2}{2\sigma^2} \right\} = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{k=1}^n (X_k - \mu)^2 \right\}.$$

Теперь будем решать пункты по отдельности.

- (a) Заметим, что

$$\sum_{k=1}^n (X_k - \mu)^2 = \sum_{k=1}^n X_k^2 - 2\mu \sum_{k=1}^n X_k + n\mu^2.$$

Тогда

$$p_\theta(\mathbf{X}) = \frac{1}{(2\pi)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{k=1}^n X_k^2 + \frac{\mu}{\sigma^2} \sum_{k=1}^n X_k - \frac{n\mu^2}{2\sigma^2} - \frac{n}{2} \ln \sigma^2 \right\}.$$

Следовательно, это экспоненциальное семейство, в котором $a_1(\mu) = \mu/(2\sigma^2)$, $u_1(\mathbf{X}) = \sum_{k=1}^n X_k$, $h(\mathbf{X}) = (2\pi)^{-n/2} \exp \{ -(2\sigma)^{-1} \sum_{k=1}^n X_k^2 \}$, а $v(\mu) = -n\mu^2/(2\sigma^2) - n \ln \sigma$. Так как θ пробегает всю прямую, то $a_1(\mu)$ тоже пробегает всю прямую, а, следовательно, и какой-то отрезок. Тогда $u_1(\mathbf{X})$ есть полная достаточная статистика. Теперь достаточно решить уравнение несмещённости, но его уже решали очень много раз — подходит выборочное среднее $\bar{\mathbf{X}}$, что и будет оптимальной оценкой среднего.

- (b) В данном случае никакие преобразования не требуются. Это будет экспоненциальное семейство, в котором $a_1(\theta) = -1/(2\sigma^2)$, $u_1(\mathbf{X}) = \sum_{k=1}^n (X_k - \mu)^2$, $v(\theta) = -n \ln \sigma$, $h(\mathbf{X}) = (2\pi)^{-n/2}$. В данном случае $a_1(\theta)$ пробегает $(-\infty, 0)$, поэтому $u_1(\mathbf{X})$ есть полная достаточная статистика. Теперь нужно решить уравнение несмещённости. Для этого заметим, что

$$\mathbb{E}_\theta \left[\sum_{k=1}^n (X_k - \mu)^2 \right] = n \mathbb{D}_\theta[X_1] = n\sigma^2.$$

Тогда оптимальной оценкой σ^2 будет $n^{-1} \sum_{k=1}^n (X_k - \mu)^2$.

(с) В данном случае запишем плотность следующим образом:

$$p_{\theta}(\mathbf{X}) = \frac{1}{(2\pi)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{k=1}^n X_k^2 + \frac{\mu}{\sigma^2} \sum_{k=1}^n X_k - \frac{n\mu^2}{2\sigma^2} - \frac{n}{2} \ln \sigma^2 \right\}.$$

Тогда это будет экспоненциальное семейство распределений с $\mathbf{a}(\theta) = (\mu/(2\sigma^2), -1/(2\sigma^2))$, $\mathbf{u}(\mathbf{X}) = (\sum_{k=1}^n X_k, \sum_{k=1}^n X_k^2)$, $v(\theta) = -n\mu^2/(2\sigma^2) - n \ln \sigma$, $h(\mathbf{X}) = (2\pi)^{-n/2}$. Заметим, что $\mathbf{a}(\theta)$ пробегает $\mathbb{R} \times \mathbb{R}_{++}$. Тогда $\mathbf{u}(\mathbf{X})$ есть полная и достаточная статистика. Теперь надо решить уравнения несмещённости. Сразу же скажем, что $E_{\theta}[\bar{X}] = \mu$, то есть выборочное среднее есть оптимальная оценка среднего. Далее,

$$E_{\theta} \left[\frac{1}{n} \sum_{k=1}^n X_k^2 \right] = \mu^2 + \sigma^2$$

Нужно избавиться от μ^2 . Для этого посчитаем матожидание квадрата выборочно среднего, пользуясь тем, что $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$:

$$E_{\theta}[(\bar{X})^2] = \mu^2 + \frac{\sigma^2}{n}.$$

Поэтому оптимальной оценкой σ^2 будет

$$S^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2. \quad \square$$

9 Лекция 9

9.1 Линейная регрессия

Сегодня мы разберём классическую задачу линейной регрессии. Обычно она формулируется следующим образом: допустим, что мы хотим узнать какую-то величину $\ell \in \mathbb{R}^n$ и проводим измерения. Но реальные инструменты не могут получить абсолютную точность, поэтому замеры будут шумными:

$$\mathbf{X} = \ell + \varepsilon,$$

где \mathbf{X} — результаты измерений, а ε — ошибка измерения. Обычно предполагается, что ε есть гауссовский вектор, но сегодня мы будем рассматривать общий случай. Далее, на ε накладываются следующие ограничения:

$$E[\varepsilon] = \mathbf{0}, \quad D[\varepsilon] = \sigma^2 \mathbf{I}_n.$$

То есть аппарат для наблюдения достаточно хорош, чтобы в среднем скос был нулевой, а ошибки между собой некоррелируют и имеют одинаковую неизвестную дисперсию σ^2 . В таком виде можно представить почти все реальные задачи.

В чём состоит линейность модели? В том, что мы кое-что знаем про ℓ : это не какой-то произвольный вектор, а вектор из нетривиального *известного* линейного подпространства L размерности $k < n$. Но что означает, что мы знаем L ? То, что нам известен его некоторый базис $\mathbf{z}_1, \dots, \mathbf{z}_k \in \mathbb{R}^n$. Составим из базисных векторов матрицу $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_k) \in \mathbb{R}^{n \times k}$. Понятно, что в таком случае $\text{rank } \mathbf{Z} = k$. Тогда ℓ можно разложить в линейную комбинацию:

$$\ell = \sum_{i=1}^k \theta_i \mathbf{z}_i = \mathbf{Z}\boldsymbol{\theta}, \quad \text{где } \boldsymbol{\theta} = (\theta_1, \dots, \theta_k).$$

Смысл вектора $\boldsymbol{\theta}$ очевиден: это неизвестные координаты вектора ℓ в базисе $\mathbf{z}_1, \dots, \mathbf{z}_k$. Тогда задача линейной регрессии переформулируется следующим образом: по наблюдению \mathbf{X} нужно оценить $\boldsymbol{\theta}$ и σ^2 .

Рассмотрим какие-нибудь примеры:

- Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из нормального распределения $\mathcal{N}(\mu, \sigma^2)$. В таком случае задачей линейной регрессии можно считать задачу оценки параметров. Для этого скажем, что

$$\mathbf{X} = \begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix} + \begin{pmatrix} X_1 - \mu \\ \vdots \\ X_n - \mu \end{pmatrix}.$$

Далее заметим, что первый вектор лежит в одномерном линейном подпространстве с базисным вектором $(1, \dots, 1)$, а второй вектор содержит независимые и одинаково распределённые случайные величины с нормальным распределением $\mathcal{N}(0, \sigma^2)$. Тем самым получаем задачу линейной регрессии.

- Теперь рассмотрим пример посложнее. Пусть тело движется по прямой с постоянной скоростью. Далее, в некоторые моменты времени t_1, \dots, t_n мы измеряем координату тела:

$$X_i = a + bt_i + \varepsilon_i,$$

где a — стартовая позиция, b — скорость движения тела и ε_i — погрешность измерения на i -й итерации. Задача состоит в оценке неизвестных a и b . Как привести эту задачу к задаче линейной регрессии? Для этого поймём, как будет выглядеть вектор ℓ :

$$\ell = \begin{pmatrix} a + bt_1 \\ \vdots \\ a + bt_n \end{pmatrix} = \begin{pmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_n \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix}$$

Тем самым разумно положить $\theta = (a, b)$: всё неизвестное в задаче, кроме σ^2 , заносится в θ . Далее, матрицу \mathbf{Z} можно ввести так:

$$\mathbf{Z} = \begin{pmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_n \end{pmatrix}.$$

Если не все наблюдения были проведены в одно и то же время, то это настоящий двумерный базис и получается легальная задача линейной регрессии. Далее, можно провести аналогичные рассуждения, если рассматривать ракету, которая летит по параболе. В таком случае будет добавляться член квадратичной зависимости от времени ct_i^2 , но суть от этого не изменится.

Как было сказано ранее, задача состоит в оценивании пары (θ, σ^2) . Будем использовать *метод наименьших квадратов*. Вообще, по этому методу можно оценивать почти всё, что угодно, но в задаче линейной регрессии он даёт содержательные свойства. Сформулируем сам метод:

Определение 29. Оценкой θ по методу наименьших квадратов называется

$$\hat{\theta}(\mathbf{X}) = \arg \min_{\theta \in \mathbb{R}^k} \|\mathbf{X} - \mathbf{Z}\theta\|^2.$$

Геометрический смысл этого определения крайне прост: $\mathbf{Z}\hat{\theta}(\mathbf{X})$ есть проекция \mathbf{X} на L . Но тогда $\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})$ есть проекция \mathbf{X} на ортогональное подпространство L^\perp . Отсюда можно получить замкнутую форму для $\hat{\theta}(\mathbf{X})$.

Лемма. Оценка наименьших квадратов имеет следующий вид:

$$\hat{\theta}(\mathbf{X}) = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}.$$

Перед этим заметим, что $\mathbf{Z}^\top \mathbf{Z} \in \mathbb{R}^{k \times k}$ является обратимой матрицей, так как столбцы \mathbf{Z} являются базисными.

Доказательство. Как было сказано выше, $\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})$ есть проекция \mathbf{X} на L^\perp . Следовательно, для любого $\theta \in \mathbb{R}^k$ $\mathbf{Z}\theta$ ортогонально $\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})$:

$$\langle \mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X}), \mathbf{Z}\theta \rangle = \theta^\top \mathbf{Z}^\top (\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})) = 0.$$

Но второе равенство выполнено для всех $\theta \in \mathbb{R}^k$ тогда и только тогда, когда $\mathbf{Z}^\top (\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})) = \mathbf{0}$, так как иначе можно взять $\theta = \mathbf{Z}^\top (\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})) \neq \mathbf{0}$ и получить $0 = \theta^\top \theta > 0$. Следовательно,

$$\mathbf{Z}^\top \mathbf{X} = \mathbf{Z}^\top \mathbf{Z}\hat{\theta}(\mathbf{X}) \implies \hat{\theta}(\mathbf{X}) = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}. \quad \square$$

Вернёмся к примерам и попробуем посчитать оценки:

- В первом случае матрица \mathbf{Z} равна $(1, 1, \dots, 1)^\top \in \mathbb{R}^{n \times 1}$. Тогда $\mathbf{Z}^\top \mathbf{Z} = n$ и

$$\hat{\mu}(\mathbf{X}) = \frac{1}{n} \mathbf{Z}^\top \mathbf{X} = \sum_{i=1}^n X_i = \bar{X}.$$

- Для начала посчитаем $\mathbf{Z}^\top \mathbf{Z}$:

$$\mathbf{Z}^\top \mathbf{Z} = \begin{pmatrix} n & \sum_{i=1}^n t_i \\ \sum_{i=1}^n t_i & \sum_{i=1}^n t_i^2 \end{pmatrix}.$$

Теперь обратим её:

$$(\mathbf{Z}^\top \mathbf{Z})^{-1} = \frac{1}{n \sum_{i=1}^n t_i^2 - (\sum_{i=1}^n t_i)^2} \begin{pmatrix} \sum_{i=1}^n t_i^2 & -\sum_{i=1}^n t_i \\ -\sum_{i=1}^n t_i & n \end{pmatrix}.$$

Далее, посчитаем $\mathbf{Z}^\top \mathbf{X}$:

$$\mathbf{Z}^\top \mathbf{X} = \begin{pmatrix} \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i t_i \end{pmatrix}$$

Осталось помножить:

$$(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X} = \frac{1}{n \sum_{i=1}^n t_i^2 - (\sum_{i=1}^n t_i)^2} \begin{pmatrix} (\sum_{i=1}^n X_i)(\sum_{i=1}^n t_i^2) - (\sum_{i=1}^n t_i)(\sum_{i=1}^n X_i t_i) \\ n \sum_{i=1}^n X_i t_i - (\sum_{i=1}^n X_i)(\sum_{i=1}^n t_i) \end{pmatrix}$$

Эта формула хороша всем, кроме одного: вычислять ответ с её помощью может оказаться непомерно дорого, так как вычисление подразумевает перемножение и обращение матриц: матрицу 2×2 обернуть достаточно быстро, а матрицу 10000×10000 уже долго. Поэтому чаще оценку наименьших квадратов ищут по определению, то есть решают задачу минимизации квадратичной функции.

Попытаемся понять, какие свойства могут быть у оценок методом наименьших квадратов. Начнём с самого простого: посчитаем матожидание и дисперсию.

Лемма. Для оценки методом наименьших квадратов $\hat{\theta}(\mathbf{X})$ выполнено следующее:

$$\mathbb{E}[\hat{\theta}(\mathbf{X})] = \theta, \quad \mathbb{D}[\hat{\theta}(\mathbf{X})] = \sigma^2 (\mathbf{Z}^\top \mathbf{Z})^{-1}.$$

Доказательство. Для начала заметим, что матрица $\mathbf{Z}^\top \mathbf{Z}$ симметрична. Далее, $\mathbf{X} = \mathbf{Z}\theta + \varepsilon$, поэтому $\mathbb{E}[\mathbf{X}] = \mathbf{Z}\theta$ и $\mathbb{D}[\mathbf{X}] = \sigma^2 \mathbf{I}_n$. Тогда

$$\begin{aligned} \mathbb{E}[\hat{\theta}(\mathbf{X})] &= \mathbb{E}[(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}] = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbb{E}[\mathbf{X}] = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Z}\theta = \theta, \\ \mathbb{D}[\hat{\theta}(\mathbf{X})] &= \mathbb{D}[(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}] = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbb{D}[\mathbf{X}] \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} = \sigma^2 (\mathbf{Z}^\top \mathbf{Z})^{-1}. \end{aligned} \quad \square$$

Можно ли что-нибудь ещё сказать про эту оценку? Асимптотические свойства рассматривать бесполезно, так как не понятно, что будет происходить с оценкой при увеличении размерности ℓ . Однако ту же оптимальность можно проверить. И оказывается, что она оптимальна, но в немного более узком классе, чем обычно.

Лемма. Оценка методом наименьших квадратов $\hat{\theta}(\mathbf{X})$ является оптимальной оценкой θ в классе линейных несмещённых оценок, то есть оценок вида $\mathbf{B}\mathbf{X}$, где \mathbf{B} — неслучайная матрица.

Доказательство. Пусть $\theta^*(\mathbf{X}) = \mathbf{B}\mathbf{X}$ — другая несмещённая оценка θ . Тогда для любого $\theta \in \mathbb{R}^k$

$$\theta = \mathbb{E}_\theta[\theta^*(\mathbf{X})] = \mathbf{B} \mathbb{E}_\theta[\mathbf{X}] = \mathbf{B}\mathbf{Z}\theta.$$

Далее, $\mathbf{B} \in \mathbb{R}^{k \times n}$, $\mathbf{Z} \in \mathbb{R}^{n \times k}$ и $\mathbf{B}\mathbf{Z} \in \mathbb{R}^{k \times k}$. Но тогда $\mathbf{B}\mathbf{Z} = \mathbf{I}_k$. Далее, заметим, что

$$\mathbb{D}_\theta[\mathbf{B}\mathbf{X}] = \mathbf{B} \mathbb{D}_\theta[\mathbf{X}] \mathbf{B}^\top = \sigma^2 \mathbf{B}\mathbf{B}^\top.$$

В итоге нужно доказать, что $\sigma^2 \mathbf{B}\mathbf{B}^\top \succcurlyeq \sigma^2 (\mathbf{Z}^\top \mathbf{Z})^{-1}$. Для этого рассмотрим следующую дисперсию:

$$\begin{aligned} \mathbb{D}_\theta[(\mathbf{B} - (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top) \mathbf{X}] &= (\mathbf{B} - (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top) \sigma^2 \mathbf{I}_n (\mathbf{B}^\top - \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1}) \\ &= \sigma^2 (\mathbf{B}\mathbf{B}^\top - \mathbf{B}\mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} - (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{B}^\top + (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1}) \\ &= \sigma^2 \mathbf{B}\mathbf{B}^\top - \sigma^2 (\mathbf{Z}^\top \mathbf{Z})^{-1} \succcurlyeq \mathbf{0}. \end{aligned} \quad \square$$

Отлично, мы получили наилучшую линейную оценку для θ . Но ещё есть неизвестный параметр σ^2 , который оценить с помощью метода наименьших квадратов не получится. И что делать? Оказывается, это можно сделать следующим образом.

Лемма. Пусть $\hat{\theta}(\mathbf{X})$ — это оценка методом наименьших квадратов. Тогда

$$\mathbb{E}_\theta[\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2] = (n - k)\sigma^2.$$

Доказательство. Для начала заметим, что $E_{\theta}[\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})] = \mathbf{0}$. Тогда

$$\begin{aligned} E_{\theta}[\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2] &= \sum_{i=1}^n D_{\theta}[X_i - (\mathbf{Z}\theta)_i] = \text{tr } D_{\theta}[\mathbf{X} - \mathbf{Z}\theta] = \text{tr } D_{\theta}[(\mathbf{I}_n - \mathbf{Z}(\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z}^{\top})\mathbf{X}] \\ &= \sigma^2 \text{tr}(\mathbf{I}_n - \mathbf{Z}(\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z}^{\top})(\mathbf{I}_n - \mathbf{Z}(\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z}^{\top}) \\ &= \sigma^2 \text{tr}(\mathbf{I}_n - 2\mathbf{Z}^{\top}(\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z} + \mathbf{Z}(\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z}^{\top}\mathbf{Z}(\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z}^{\top}) \\ &= \sigma^2 \text{tr}(\mathbf{I}_n - \mathbf{Z}^{\top}(\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z}) = \sigma^2(n - \text{tr}((\mathbf{Z}^{\top}\mathbf{Z})^{-1}\mathbf{Z}^{\top}\mathbf{Z})) = \sigma^2(n - k). \end{aligned} \quad \square$$

Тем самым получаем, что несмещённой оценкой σ^2 является

$$\hat{\sigma}^2(\mathbf{X}) = \frac{1}{n - k} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2.$$

Напоследок рассмотрим простую задачу на линейную регрессию.

Задача 24. Допустим, что есть два груза с массами a и b . Их взвесили по отдельности и вместе. Результаты записали в вектор (X_1, X_2, X_3) . Найти оценку для (a, b) методом наименьших квадратов.

Решение. Для начала сведём это к задаче линейной регрессии:

$$\underbrace{\begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}}_{\mathbf{X}} = \begin{pmatrix} a + \varepsilon_1 \\ b + \varepsilon_2 \\ a + b + \varepsilon_3 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}}_{\mathbf{Z}} \underbrace{\begin{pmatrix} a \\ b \end{pmatrix}}_{\boldsymbol{\theta}} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \end{pmatrix}}_{\boldsymbol{\varepsilon}}.$$

Теперь можно считать по уже известной формуле. Для начала посчитаем $(\mathbf{Z}^{\top}\mathbf{Z})^{-1}$:

$$\mathbf{Z}^{\top}\mathbf{Z} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \implies (\mathbf{Z}^{\top}\mathbf{Z})^{-1} = \frac{1}{3} \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}.$$

Далее,

$$\mathbf{Z}^{\top}\mathbf{X} = \begin{pmatrix} X_1 + X_3 \\ X_2 + X_3 \end{pmatrix}$$

Следовательно,

$$\hat{\theta}(\mathbf{X}) = \frac{1}{3} \begin{pmatrix} 2X_1 - X_2 + X_3 \\ -X_1 + 2X_2 + X_3 \end{pmatrix} \implies \begin{cases} \hat{a}(\mathbf{X}) = (2X_1 - X_2 + X_3)/3 \\ \hat{b}(\mathbf{X}) = (-X_1 + 2X_2 + X_3)/3 \end{cases}$$

□

9.2 Гауссовская линейная регрессия

Далее в общем случае какие-либо предположения делать уже крайне сложно, поэтому начнём добавлять какие-то ограничения на $\boldsymbol{\varepsilon}$. Самое популярное — затребовать нормальность распределения $\boldsymbol{\varepsilon}$. В таком случае линейную регрессию называют гауссовской. В ней

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n) \implies \mathbf{X} \sim \mathcal{N}(\boldsymbol{\ell}, \sigma^2 \mathbf{I}_n) = \mathcal{N}(\mathbf{Z}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_n).$$

Оказывается, что в таком случае полученная ранее оценка методом наименьших квадратов обладает крайне приятными свойствами.

Теорема 24. В гауссовской линейной регрессионной модели пара $(\hat{\theta}(\mathbf{X}), \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2)$ является полной достаточной статистикой для $(\boldsymbol{\theta}, \sigma^2)$.

Доказательство. Найдём плотность \mathbf{X} :

$$\begin{aligned} p_{\mathbf{X}}(\mathbf{x}) &= \frac{1}{(2\pi)^{n/2} \sqrt{|\sigma^2 \mathbf{I}_n|}} \exp \left\{ -\frac{1}{2} (\mathbf{X} - \mathbf{Z}\boldsymbol{\theta})^{\top} (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}) \right\} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2 \right\} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X}) + \mathbf{Z}\hat{\theta}(\mathbf{X}) - \mathbf{Z}\boldsymbol{\theta}\|^2 \right\}. \end{aligned}$$

Далее нужно приглядеться к этой норме. Вспомним, что $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть проекция \mathbf{X} на L^\perp , а $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X}) - \mathbf{Z}\boldsymbol{\theta}$ лежит в L . Но тогда по теореме Пифагора

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2}{2\sigma^2} - \frac{\|\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X}) - \mathbf{Z}\boldsymbol{\theta}\|^2}{\sigma^2} \right\}.$$

Отсюда уже видно то, что по критерию факторизации $(\hat{\boldsymbol{\theta}}(\mathbf{X}), \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2)$ есть достаточная статистика. Но нам нужна полнота, поэтому нужно несколько преобразовать плотность и привести её к виду экспоненциального семейства. Для этого немного откатимся назад и раскроем:

$$\begin{aligned} p_{\mathbf{X}}(\mathbf{x}) &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2 \right\} = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{\|\mathbf{X}\|^2}{2\sigma^2} + \frac{\langle \mathbf{Z}^\top \mathbf{X}, \boldsymbol{\theta} \rangle}{\sigma^2} - \frac{\|\mathbf{Z}\boldsymbol{\theta}\|^2}{2\sigma^2} \right\} \\ &= \frac{1}{(2\pi)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{X}\|^2 + \sum_{i=1}^k (\mathbf{Z}^\top \mathbf{X})_i \frac{\theta_i}{\sigma^2} - \frac{\|\mathbf{Z}\boldsymbol{\theta}\|^2}{2\sigma^2} - \frac{n}{2} \ln \sigma^2 \right\}. \end{aligned}$$

Заметим, что при пробегании $\boldsymbol{\theta} \in \mathbb{R}^k$ и $\sigma^2 > 0$ вектор $(-1/(2\sigma^2), \theta_1/\sigma^2, \dots, \theta_k/\sigma^2)$ заметает подпространство \mathbb{R}^{k+1} . Тем самым по теореме об экспоненциальном семействе статистика $(\mathbf{Z}^\top \mathbf{X}, \|\mathbf{X}\|^2)$ является полной и достаточной статистикой для пары $(\boldsymbol{\theta}, \sigma^2)$. Осталось доказать, что между этой статистикой и статистикой из условия теоремы есть биекция. Это мы докажем на следующей лекции. \square

10 Лекция 10

Продолжение доказательства. Построим биекцию между $(\mathbf{Z}^\top \mathbf{X}, \|\mathbf{X}\|^2)$ и $(\hat{\boldsymbol{\theta}}(\mathbf{X}), \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2)$. Заметим, что по теореме Пифагора

$$\|\mathbf{X}\|^2 = \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 + \|\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2,$$

так как $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ лежит в L^\perp , а $\|\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2$ — в L . Далее, согласно формуле для оценки методом наименьших квадратов

$$\mathbf{Z}^\top \mathbf{X} = (\mathbf{Z}^\top \mathbf{Z}) \hat{\boldsymbol{\theta}}(\mathbf{X}).$$

Тем самым, если мы знаем значения пары $(\hat{\boldsymbol{\theta}}(\mathbf{X}), \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2)$, то мы знаем значения пары $(\mathbf{Z}^\top \mathbf{X}, \|\mathbf{X}\|^2)$ (так как \mathbf{Z} тоже известна). В обратную сторону всё аналогично:

$$\begin{aligned} \hat{\boldsymbol{\theta}}(\mathbf{X}) &= (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}, \\ \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 &= \|\mathbf{X}\|^2 - \|\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2. \end{aligned}$$

Тем самым получаем полную взаимнообратную связь между парами. Следовательно, $(\hat{\boldsymbol{\theta}}(\mathbf{X}), \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2)$ есть полная достаточная статистика. \square

Теперь рассмотрим некоторые следствия данной теоремы.

Следствие. $\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть оптимальная оценка для $\boldsymbol{\theta}$, $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ — для $\boldsymbol{\ell}$, а $\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2/(n-k)$ — для σ^2 .

Доказательство. Ранее было показано, что эти функции являются несмещёнными оценками для соответствующих параметров. Но они все являются функциями от полной и достаточной статистики. Следовательно, они оптимальны. \square

10.1 Доверительные интервалы и области

Теперь мы бы хотели построить доверительные интервалы для параметров гауссовской линейной регрессии. Но для этого нужна теорема об ортогональных разложениях. Докажем её, но перед этим вспомним, что такое распределение хи-квадрат.

Определение 30. Будем говорить, что случайная величина ξ имеет *распределение хи-квадрат с n степенями свободы*, если её плотность равна

$$p_\xi(x) = \frac{x^{n/2-1} e^{-x/2}}{2^{n/2} \Gamma(n/2)} \mathbf{1}_{x \geq 0},$$

то есть $\xi \sim \Gamma(n/2, 1/2)$. Обозначение: $\xi \sim \chi_n^2$.

Далее вспомним ещё один простой факт: если ξ_1, \dots, ξ_n есть независимые и одинаково распределённые случайные величины с стандартным нормальным распределением, то $\xi_1^2 + \dots + \xi_n^2 \sim \chi_n^2$. Теперь можно приступить к теореме.

Теорема 25 (об ортогональных разложениях). Пусть $\mathbf{X} \sim \mathcal{N}(\ell, \sigma^2 \mathbf{I}_n)$, а $L_1 \oplus L_2 \oplus \dots \oplus L_r$ есть разложение \mathbb{R}^n в прямую сумму линейных ортогональных подпространств. Далее, пусть \mathbf{Y}_i есть проекция \mathbf{X} на L_i . Тогда Y_1, \dots, Y_r независимы в совокупности и

$$\frac{1}{\sigma^2} \|\mathbf{Y}_i - \mathbb{E}[\mathbf{Y}_i]\|^2 \sim \chi_{\dim L_i}^2.$$

Доказательство. Идея состоит в том, чтобы свести эту теорему к факту, сформулированному выше. Для этого возьмём в \mathbb{R}^n ортонормированный базис $\mathbf{f}_1, \dots, \mathbf{f}_n$ такой, что

$$\begin{aligned} \mathbf{f}_1, \dots, \mathbf{f}_{k_1} &\text{— базис } L_1, \\ \mathbf{f}_{k_1+1}, \dots, \mathbf{f}_{k_1+k_2} &\text{— базис } L_2, \\ &\vdots \\ \mathbf{f}_{k_1+\dots+k_{r-1}+1}, \dots, \mathbf{f}_n &\text{— базис } L_r. \end{aligned}$$

В таком случае размерность L_i равна k_i для всех $i = 1, \dots, r$. Далее, введём коэффициенты $W_i = \langle \mathbf{X}, \mathbf{f}_i \rangle$ для всех $i = 1, \dots, n$. Скомпонуем эти коэффициенты в вектор и представим его в виде линейного преобразования \mathbf{X} :

$$\mathbf{W} = \begin{pmatrix} W_1 \\ \vdots \\ W_n \end{pmatrix} = \mathbf{C}\mathbf{X}, \text{ где } \mathbf{C} = \begin{pmatrix} \mathbf{f}_1^\top \\ \vdots \\ \mathbf{f}_n^\top \end{pmatrix}.$$

Заметим, что матрица \mathbf{C} ортогональна, так как её строки задают ортонормированный базис \mathbb{R}^n . Так как \mathbf{X} гауссовский, то и \mathbf{W} тоже будет гауссовским, так как он получается линейным преобразованием. Осталось заметить, что $\mathbf{W} \sim \mathcal{N}(\mathbf{C}\ell, \sigma^2 \mathbf{I}_n)$, так как $\mathbf{C}\mathbf{C}^\top = \mathbf{I}_n$. Тем самым компоненты \mathbf{W} независимы в совокупности. Тогда \mathbf{Y}_i , которые считаются следующим образом:

$$\mathbf{Y}_j = W_{k_1+\dots+k_{j-1}+1} \mathbf{f}_{k_1+\dots+k_{j-1}+1} + \dots + W_{k_1+\dots+k_j} \mathbf{f}_{k_1+\dots+k_j},$$

тоже будут независимы в совокупности. Теперь докажем свойство про их распределение. Для этого заметим, что для всех $i = 1, \dots, n$

$$\frac{W_i - \mathbb{E}[W_i]}{\sigma} \sim \mathcal{N}(0, 1).$$

В таком случае, так как векторы ортонормированные, то

$$\begin{aligned} \frac{1}{\sigma^2} \|\mathbf{Y}_j - \mathbb{E}[\mathbf{Y}_j]\|^2 &= \left\| \sum_{i=1}^{k_j} \frac{W_{k_1+\dots+k_{j-1}+i} - \mathbb{E}[W_{k_1+\dots+k_{j-1}+i}]}{\sigma} \mathbf{f}_{k_1+\dots+k_{j-1}+i} \right\|^2 \\ &= \sum_{i=1}^{k_j} \left(\frac{W_{k_1+\dots+k_{j-1}+i} - \mathbb{E}[W_{k_1+\dots+k_{j-1}+i}]}{\sigma} \right)^2 \sim \chi_{k_j}^2. \end{aligned} \quad \square$$

Теперь выпишем из этого достаточно важное следствие:

Следствие. В гауссовской линейной модели $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ (и $\hat{\boldsymbol{\theta}}(\mathbf{X})$) независимо с $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ и

$$\begin{aligned} \frac{1}{\sigma^2} \|\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X}) - \mathbf{Z}\boldsymbol{\theta}\|^2 &\sim \chi_k^2, \\ \frac{1}{\sigma^2} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 &\sim \chi_{n-k}^2. \end{aligned}$$

Доказательство. Как известно, $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть проекция \mathbf{X} на L^\perp , а $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ проецирует \mathbf{X} на L . Так как $\mathbb{R}^n = L \oplus L^\perp$, то применима теорема об ортогональных разложениях. Следовательно, $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ и $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ независимы и верны выражения из условия. \square

Теперь приступим к доверительным интервалам. Перед этим надо определить, что это вообще такое.

Определение 31. Пусть $\{P_\theta \mid \theta \in \Theta\}$ — параметрическое семейство, \mathbf{X} — наблюдение из параметрического распределения $P \in \{P_\theta \mid \theta \in \Theta\}$, и параметр θ одномерный: $\Theta \subseteq \mathbb{R}$. Будем называть пару статистик $(T_1(\mathbf{X}), T_2(\mathbf{X}))$ *доверительным интервалом* для θ с уровнем доверия γ , если для всех $\theta \in \Theta$

$$P_\theta(T_1(\mathbf{X}) < \theta < T_2(\mathbf{X})) \geq \gamma.$$

Если в данном неравенстве всегда достигается равенство, то доверительный интервал называют *точным*.

Данное определение почти без изменений переносится на многомерный параметр.

Определение 32. Подмножество $S(\mathbf{X}) \subseteq \Theta$ называется доверительной областью для θ с уровнем доверия γ , если для всех $\theta \in \Theta$

$$P_\theta(\theta \in S(\mathbf{X})) \geq \gamma.$$

Ну что же, сперва построим доверительный интервал для σ^2 . Для этого нужно воспользоваться так называемой центральной статистикой. Дадим определение:

Определение 33. Функция от выборки $S_\theta(\mathbf{X})$, которая может зависеть от параметра θ , называется *центральной статистикой*, если её распределение не зависит от θ .

Заметим, что в данном случае можно использовать $\sigma^{-2} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2$ в качестве центральной статистики, так как ранее было доказано, что она имеет распределение хи-квадрат с $n - k$ степенями свободы. Пусть u_α равно α -квантилю χ^2_{n-k} . Тогда

$$P\left(u_{(1-\gamma)/2} < \frac{1}{\sigma^2} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2 < u_{(1+\gamma)/2}\right) = \gamma.$$

Следовательно, точный доверительный интервал для σ^2 уровня доверия γ имеет вид

$$\left[\frac{1}{u_{(1+\gamma)/2}} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2, \frac{1}{u_{(1-\gamma)/2}} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2 \right]$$

Теперь попытаемся определить доверительный интервал для θ_i . Для этого заметим, что $\hat{\theta}(\mathbf{X}) \sim \mathcal{N}(\theta, \sigma^2 \mathbf{A})$, где $\mathbf{A} = (\mathbf{Z}^\top \mathbf{Z})^{-1}$. В таком случае

$$\frac{\hat{\theta}_i - \theta_i}{\sqrt{\sigma^2 a_{ii}}} \sim \mathcal{N}(0, 1).$$

Но возникает проблема: всё ещё остаётся σ^2 . Что делать? Заменить на оценку. Оказывается, что в гауссовском случае всё остаётся весьма хорошо. Полученная случайная величина будет иметь распределение Стьюдента с $n - k$ степенями свободы:

$$\sqrt{\frac{n-k}{a_{ii}}} \frac{\hat{\theta}_i - \theta_i}{\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|} \sim T_{n-k}.$$

Почему это так? Это следует из того, что $\hat{\theta}_i$ не зависит от $\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|$, и дробь имеет вид отношения стандартной нормальной случайной величины к корню из случайной величины с распределением хи-квадрат, что и есть распределение Стьюдента по определению.

Пусть t_α есть α -квантиль распределения Стьюдента T_{n-k} . Так как оно симметрично относительно нуля, то можно сразу записать точный доверительный интервал с уровнем доверия γ в следующем виде:

$$\hat{\theta}_i - t_{(1+\gamma)/2} \sqrt{\frac{a_{ii}}{n-k}} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\| < \theta_i < \hat{\theta}_i + t_{(1+\gamma)/2} \sqrt{\frac{a_{ii}}{n-k}} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|.$$

Осталось построить доверительную область для θ . Здесь возникает третье стандартное распределение, связанное с гауссовским — распределение Фишера.

Определение 34. Пусть ξ, η — независимые случайные величины, причём $\xi \sim \chi^2_k$, $\eta \sim \chi^2_n$. Тогда $\delta = \frac{\xi/k}{\eta/n}$ имеет распределение Фишера с (k, n) степенями свободы. Обозначение: $\delta \sim F_{k,n}$.

Согласно следствию из теоремы об ортогональных разложениях

$$\frac{\|\mathbf{Z}\hat{\theta}(\mathbf{X}) - \mathbf{Z}\theta\|^2}{\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2} \frac{n-k}{k} \sim F_{k,n-k}.$$

Далее, пусть f_α есть α -квантиль распределения $F_{k,n-k}$. Так как оно неотрицательно, то доверительной областью будет

$$S(\mathbf{X}) = \left\{ \theta \in \mathbb{R}^k : \frac{\|\mathbf{Z}\hat{\theta}(\mathbf{X}) - \mathbf{Z}\theta\|^2}{\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2} \frac{n-k}{k} < f_\gamma \right\}.$$

Стоит заметить, что геометрически доверительная область будет иметь форму эллипсоида.

10.2 Проверка гипотез

На этом мы пока что простимся с линейной регрессией и вернёмся к ней несколько позже. Пока что разберёмся с такой вещью, как проверка гипотез. Вообще, проверка гипотез — это больше методология, тут математики меньше и важнее понять сам метод. Что вообще предполагается в статистике? Предполагается следующее: если вероятность события крайне мала (меньше, чем 0.01, например), то в единичном испытании данное событие не происходит никогда. Если мы так не считаем, то могут происходить события крайне малой вероятности и наша статистика нужно выбросить и думать дальше. Примерно на этом построена методология проверки гипотез: что мы должны с чем-то согласиться и действовать дальше. Обычно согласуют некоторый уровень ошибки.

Теперь формализуем вышесказанное. Пусть \mathbf{X} — наблюдение с неизвестным распределением $P \in \mathcal{P}$, где \mathcal{P} есть некоторое семейство распределений. Далее, возьмём подсемейство $\mathcal{P}_0 \subset \mathcal{P}$. Поставим следующий вопрос: верно ли то, что $P \in \mathcal{P}_0$? Ответ на этот вопрос позволяет сузить семейство распределений, что упростит анализ. Если ответ положительный, то \mathcal{P} можно сузить до \mathcal{P}_0 . Если же ответ отрицательный, то \mathcal{P} сужается до $\mathcal{P} \setminus \mathcal{P}_0$. Компромисс в каком-то смысле невозможен: для получения компромисса нужно набрать ещё данных. Рассуждение за этим следующее: по текущим данным оказалось, что верна одна гипотеза, хотя кажется, что должна быть верна другая. Кажется, что выборка оказалась плохой и нужно набрать ещё данных для повторной проверки.

Обычно статистическая гипотеза (или предположение) обозначается следующим образом:

$$H_0: P \in \mathcal{P}_0.$$

Что бывает ещё? В принципе, мы можем иметь не только одну гипотезу. Тогда параллельно с ней вводят альтернативу $H_1: P \in \mathcal{P}_1 \subseteq \mathcal{P} \setminus \mathcal{P}_0$. Альтернатив может быть несколько:

$$\begin{aligned} H_0: P \in \mathcal{P}_0, \\ H_1: P \in \mathcal{P}_1, \\ \vdots \\ H_k: P \in \mathcal{P}_k. \end{aligned}$$

Методология метода проверки гипотез следующая. Допустим, что сравнивается пара гипотез H_0 против H_1 . Если гипотеза H_0 принимается, то есть на вопрос «Верно ли то, что $P \in \mathcal{P}_0$?» даётся положительный ответ, то семейство распределений сужается до \mathcal{P}_0 . Однако, если гипотезу H_0 отвергают и задана альтернатива, то семейство сужают до \mathcal{P}_1 . Если же альтернатив много, то вариантов постановки вопроса несколько. Первый вариант состоит в том, что все гипотезы предполагаются равноправными и потом каким-то образом выбирается одна из них. Обычно же принято действовать иначе: гипотезы перебираются по очереди, и, если гипотеза H_i принимается, то семейство сужается до \mathcal{P}_i , иначе же рассматривается гипотеза H_{i+1} .

Далее возникает вопрос: какие могут быть ошибки? Например, мы можем промахнуться: например, гипотезу H_0 приняли, хотя на самом деле она не верна. Такие ошибки разделяют на ошибки первого и второго рода.

Определение 35. Пусть сравниваются гипотезы H_0 против H_1 . *Ошибкой первого рода* называется ситуация, когда H_0 отвергли, но она была верна. Аналогично, *ошибкой второго рода* называют ситуацию, когда H_0 приняли, но она была не верна.

И вопрос: какая из этих ошибок опаснее? Методологически считается, что ошибка первого рода опаснее. Предположим, что мы неверно отвергли H_0 . Тогда мы вынуждены рассматривать H_1 , которая заведомо неверна. В зависимости от результата проверки H_1 будет либо ошибка, либо дополнительная работа в виде проверки H_2 и оба варианта не вдохновляют. Тем самым будет проводиться много бесполезной работы. В случае ошибки второго рода мы принимаем H_0 , но она неверна. В таком случае вердикт может измениться, если добавить данных, и мы сможем двинуться в сторону чего-то более правильного. Другими словами, здесь шансы совершить много лишней работы сильно меньше.

Что это означает? То, что ошибку первого рода нужно как-то заведомо избежать, то есть вероятность того, что она произойдёт, должна быть достаточно мала (не больше 0.01, например). Но для оценки вероятности ошибки нужно понимать, как мы будем проверять гипотезы.

Определение 36. Пусть \mathcal{X} — выборочное пространство (множество всех возможных значений наблюдений). Тогда подмножество $S \subset \mathcal{X}$ называется *критерием* или *критическим множеством* для проверки гипотезы H_0 , если правило принятия H_0 выглядит так: H_0 отвергается тогда и только тогда, когда $\mathbf{X} \in S$.

По сути, S есть «плохое» множество, и, если \mathbf{X} попал туда, то H_0 считается неправильной и выбирается альтернатива. Как только появился критерий S , можно ввести вероятности ошибок первого и второго рода. Но для этого нужна функция мощности.

Определение 37. *Функцией мощности критерия S называется*

$$\beta(Q, S) = Q(X \in S), Q \in \mathcal{P}.$$

Далее возникает один момент: вероятность ошибки первого рода не является вероятностью, а является пачкой вероятностей. Аналогично со вторым родом.

Определение 38. Вероятностями ошибок первого рода называется набор $\{\beta(Q, S) : Q \in \mathcal{P}_0\}$.

Определение 39. Вероятностями ошибок второго рода называется набор $\{1 - \beta(Q, S) : Q \notin \mathcal{P}_0\}$.

Последний вопрос: какой критерий S будем подходить? Мы хотели построить очень хороший критерий в том смысле, что у него маленькие вероятности ошибок первого и второго рода. Это означает, что внутри критического множества функция мощности должна быть мала, а вне него должна быть велика. В идеале хочется получить нули, но это задача безнадёжна, и нужно на что-то соглашаться. Но на что? И как сказать, что один критерий лучше другого?

Для этого будем полагать, что вероятности ошибок первого рода ограничены сверху каким-то ε , который называют *уровнем значимости*. Далее, если мы зафиксировали уровень значимости, то среди всех подходящих критериев ищется тот, у которого минимальные вероятности ошибок второго рода. То, насколько получится минимизировать ошибку второго рода, сильно зависит от задачи. Хорошо себя показывают в этой истории так называемые *асимптотические критерии*: в них полагается, что у выборки растущий размер. С ними получается добиться следующего: вероятность ошибки первого рода не больше, чем ε , а вероятность второго рода стремится к нулю.

11 Лекция 11

Вернёмся к нашим гипотезам. Пусть проверяется гипотеза $H_0 : P \in \mathcal{P}_0$ против гипотезы $H_1 : P \in \mathcal{P}_1$, и есть какой-то критерий $S \subset \mathcal{X}$. Ранее мы вводили понятие уровня значимости. Формализуем его:

Определение 40. Критерий S имеет уровень значимости γ , если для любого $Q \in \mathcal{P}_0$ $\beta(Q, S) \leq \gamma$.

По сути, если γ является подходящим уровнем значимости, то и любое число, большее γ , тоже подходит. Поэтому наряду с уровнем значимости вводится минимальный уровень значимости:

Определение 41. Минимальным уровнем значимости или же размером критерия называют $\alpha(S) = \sup_{Q \in \mathcal{P}_0} \beta(Q, S)$.

Теперь можно ввести свойства критериев. Они называются так же, как и у оценок: несмещённость и состоятельность, но смысл у них несколько отличается. Введём их.

Определение 42. Критерий S для проверки гипотезы $H_0 : P \in \mathcal{P}_0$ против гипотезы $H_1 : P \in \mathcal{P}_1$ называется *несмещённым*,¹⁴ если

$$\sup_{Q \in \mathcal{P}_0} \beta(Q, S) \leq \inf_{Q \in \mathcal{P}_1} \beta(Q, S). \quad (1)$$

Данное свойство достаточно естественно: мы хотим, чтобы на \mathcal{P}_0 функция мощности была маленькой, а на \mathcal{P}_1 — большой. Поэтому разумно попросить, чтобы значения были разделены.

Определение 43. Пусть $\{S_n\}_{n=1}^\infty$ — последовательность¹⁵ критериев для проверки гипотезы $H_0 : P \in \mathcal{P}_0$ против гипотезы $H_1 : P \in \mathcal{P}_1$. Она называется *состоятельной* (или же говорят, что S_n — состоятельный критерий), если для любого $Q \in \mathcal{P}_1$ $\beta(Q, S_n) \rightarrow 1$ при $n \rightarrow \infty$.

По сути, это означает, что вероятности ошибки второго рода стремятся к нулю, что есть хорошо. Стоит сказать, что если асимптотический критерий не состоятелен, то он обычно бессмысленен. Но за состоятельность нужно платить, о чём расскажем позднее.

11.1 Сравнение критериев

Пусть есть два критерия S_1 и S_2 . Как их можно сравнивать? Сразу же стоит сказать, что их имеет смысл сравнивать только в том случае, когда их уровни значимости¹⁶ равны. Если это не так, то лучше взять критерий с меньшим уровнем значимости — ибо меньшая вероятность ошибки первого рода важнее, чем меньшая вероятность ошибки второго рода. Если же они равны, то сравниваем их по функции мощности.

¹⁴Не самое удачное название свойства, на самом деле — не понятно, что подразумевать под смещённостью.

¹⁵Стоит заметить, что S_n могут меняться и они живут в разных пространствах: $S_n = S_n(X_1, \dots, X_n) \subset \mathbb{R}^n$.

¹⁶Стоит сказать, что размеры критерия брать не стоит, так как чем меньше уровень значимости, тем выше вероятность ошибки второго рода, и нужно найти некоторый баланс. И, что самое интересное, его можно найти.

Определение 44. Пусть S_1 и S_2 — критерии уровня значимости ε для проверки гипотезы $H_0: P \in \mathcal{P}_0$ против гипотезы $H_1: P \in \mathcal{P}_1$. Будем говорить, что S_1 (равномерное) мощнее, чем S_2 , если для любого $Q \in \mathcal{P}_1$ $\beta(Q, S_1) \geq \beta(Q, S_2)$, то есть вероятность ошибки второго рода для критерия S_1 равномерно меньше, чем для S_2 .

Если же критерий оказывается мощнее остальных, то мы его будем называть наиболее мощным.

Определение 45. Критерий S для проверки гипотезы $H_0: P \in \mathcal{P}_0$ против гипотезы $H_1: P \in \mathcal{P}_1$ называется *равномерно наиболее мощным* критерием уровня значимости ε , если для любого другого критерия R уровня значимости ε S мощнее, чем R .

Теперь попытаемся понять, как искать равномерно наиболее мощные критерии уровня значимости ε . Самый простой вариант состоит в том, что рассматриваются только *простые* гипотезы, то есть гипотезы, в которых предлагаемые семейства содержат только одно распределение.

Рассмотрим проверку гипотезы $H_0: P = P_0$ против гипотезы $H_1: P = P_1$. Далее, предположим, что P_0 и P_1 имеют плотности $p_0(x)$ и $p_1(x)$ по одной и той же мере μ . По сути, оба распределения приходят из доминируемого семейства: например, различить пуассоновское и экспоненциальное распределение крайне просто — считаем, что экспоненциальное, если получили нецелое число, и пуассоновское иначе. Такой критерий никогда не ошибается.

Пусть $\lambda > 0$. Введём следующее множество $S_\lambda = \{x: p_1(x) - \lambda p_0(x) \geq 0\}$ и скажем, что это критерий. Логика за этим критерием понятна: мы отвергаем H_0 , если $p_1(x) > \lambda p_0(x)$. Это неравенство можно понимать в том смысле, что p_1 более правдоподобна. Теперь докажем одно интересное свойство, связанное с этим критерием.

Лемма (Нейман, Пирсон). Пусть критерий R таков, что $P_0(\mathbf{X} \in R) \leq P_0(\mathbf{X} \in S_\lambda)$. Тогда

$$a) P_1(\mathbf{X} \in R) \leq P_1(\mathbf{X} \in S_\lambda),$$

$$b) P_0(\mathbf{X} \in S_\lambda) \leq P_1(\mathbf{X} \in S_\lambda).$$

Доказательство. Для начала заметим, что для любого x ¹⁷

$$\mathbf{1}_{x \in R}(p_1(x) - \lambda p_0(x)) \leq \mathbf{1}_{x \in S_\lambda}(p_1(x) - \lambda p_0(x)). \quad (2)$$

Проинтегрируем неравенства по мере μ по всему выборочному пространству:

$$\int_R (p_1(x) - \lambda p_0(x)) \mu(dx) \leq \int_{S_\lambda} (p_1(x) - \lambda p_0(x)) \mu(dx), \quad (3)$$

$$P_1(\mathbf{X} \in R) - \lambda P_0(\mathbf{X} \in R) \leq P_1(\mathbf{X} \in S_\lambda) - \lambda P_0(\mathbf{X} \in S_\lambda), \quad (4)$$

$$P_1(\mathbf{X} \in R) - P_1(\mathbf{X} \in S_\lambda) \leq \lambda(P_0(\mathbf{X} \in R) - P_0(\mathbf{X} \in S_\lambda)) \leq 0. \quad (5)$$

Для второго пункта нужно рассмотреть два случая:

- Пусть $\lambda \geq 1$. Тогда для $x \in S_\lambda$ выполнено следующее неравенство: $p_1(x) \geq \lambda p_0(x) \geq p_0(x)$. Следовательно, $P_0(\mathbf{X} \in S_\lambda) \leq P_1(\mathbf{X} \in S_\lambda)$.
- Пусть $\lambda \in (0, 1)$. Тогда для $x \notin S_\lambda$ $p_1(x) < \lambda p_0(x) < p_0(x)$. Следовательно, $P_1(\mathbf{X} \in \overline{S_\lambda}) \leq P_0(\mathbf{X} \in \overline{S_\lambda})$ и $P_0(\mathbf{X} \in S_\lambda) \leq P_1(\mathbf{X} \in S_\lambda)$.

Тем самым для любого $\lambda > 0$ $P_0(\mathbf{X} \in S_\lambda) \leq P_1(\mathbf{X} \in S_\lambda)$. □

У данной леммы есть два приятных следствия:

Следствие. Если $\lambda > 0$ удовлетворяет уравнению $P_0(\mathbf{X} \in S_\lambda) = \varepsilon$, то S_λ является равномерно наиболее мощным критерием уровня значимости ε для проверки гипотезы $H_0: P = P_0$ против гипотезы $H_1: P = P_1$.

Следствие. S_λ — несмещённый критерий.

Определение 46. Пусть $\{P_\theta \mid \theta \in \Theta\}$, $\Theta \subseteq \mathbb{R}$ — это доминируемое семейство распределений с одномерным параметром и плотностью $p_\theta(x)$ по мере μ . Говорят, что это семейство имеет *монотонное отношение правдоподобия по статистике* $T(\mathbf{X})$, если для любых $\theta_1 < \theta_2$, $\theta_1, \theta_2 \in \Theta$

$$\frac{p_{\theta_2}(\mathbf{X})}{p_{\theta_1}(\mathbf{X})} = \psi_{\theta_1, \theta_2}(T(\mathbf{X})), \quad (6)$$

где $\psi_{\theta_1, \theta_2}$ — это монотонная функция и её вид монотонности одинаков для всех $\theta_1 < \theta_2$.

¹⁷ Данное рассуждение уже было в доказательстве теоремы Бахадура.

Лемма. Пусть $\psi_{\theta_1, \theta_2}$ всегда неубывает по $T(\mathbf{X})$. Тогда для всех $c \in \mathbb{R}$ и для всех $\theta_1 < \theta_2$, $\theta_1, \theta_2 \in \Theta$ выполнено, что

$$P_{\theta_1}(T(\mathbf{X}) \geq c) \leq P_{\theta_2}(T(\mathbf{X}) \geq c). \quad (7)$$

Доказательство. Считаем, что $\psi_{\theta_1, \theta_2} \geq 0$ и не убывает на \mathbb{R} . Введём множество $D = \{\mathbf{x}: T(\mathbf{x}) \geq c\}$, где c такое, что $\psi_{\theta_1, \theta_2}(c) \geq 1$. Тогда для любого $\mathbf{x} \in D$

$$p_{\theta_2}(\mathbf{x}) = \psi_{\theta_1, \theta_2}(T(\mathbf{x}))p_{\theta_1}(\mathbf{x}) \geq \psi_{\theta_1, \theta_2}(c)p_{\theta_1}(\mathbf{x}) \geq p_{\theta_1}(\mathbf{x}). \quad (8)$$

Интегрируем это по мере μ по множеству D :

$$\int_D p_{\theta_2}(\mathbf{x})\mu(d\mathbf{x}) \geq \int_D p_{\theta_1}(\mathbf{x})\mu(d\mathbf{x}) \quad (9)$$

$$P_{\theta_2}(T(\mathbf{X}) \geq c) \geq P_{\theta_1}(T(\mathbf{X}) \geq c). \quad (10)$$

Если же брать c такое, что $\psi_{\theta_1, \theta_2}(c) \in [0, 1]$, то рассматриваем множество \bar{D} и получаем то же самое. \square

Теорема 26 (о монотонном отношении правдоподобий). Допустим, что проверяется гипотеза $H_0: \theta \leq \theta_0$ против гипотезы $H_1: \theta > \theta_0$. Если $c \in \mathbb{R}$ удовлетворяет соотношению $P_{\theta_0}(T(\mathbf{X}) \geq c) = \gamma$, то критерий $S = \{\mathbf{x}: T(\mathbf{x}) \geq c\}$ является равномерно наиболее мощным критерием уровня доверия γ для проверки гипотезы H_0 против гипотезы H_1 .

Доказательство. Для начала возьмём какое-нибудь $\theta < \theta_0$. Тогда по лемме

$$P_{\theta}(T(\mathbf{X}) \geq c) \leq P_{\theta_0}(T(\mathbf{X}) \geq c) = \gamma. \quad (11)$$

Из этого следует, что критерий S имеет уровень значимости γ . Проверим то, что он равномерно наиболее мощный. Пусть R — любой другой критерий с уровнем мощности γ . Возьмём какое-нибудь $\theta_1 > \theta_0$. Нужно доказать, что

$$P_{\theta_1}(\mathbf{X} \in R) \leq P_{\theta_1}(\mathbf{X} \in S). \quad (12)$$

Согласно лемме Неймана-Пирсона равномерно наиболее мощный критерий уровня значимости γ для проверки гипотезы $H_0: \theta = \theta_0$ против гипотезы $H_1: \theta = \theta_1$ имеет вид

$$S_{\lambda} = \{\mathbf{x}: p_{\theta_1}(\mathbf{x}) - \lambda p_{\theta_0}(\mathbf{x}) \geq 0\}, \text{ где } P_{\theta_0}(\mathbf{X} \in S_{\lambda}) = \gamma. \quad (13)$$

Однако

$$S_{\lambda} = \left\{ \mathbf{x}: \frac{p_{\theta_1}(\mathbf{x})}{p_{\theta_0}(\mathbf{x})} \geq \lambda \right\} = \{\mathbf{x}: \psi_{\theta_1, \theta_0}(T(\mathbf{X})) \geq \lambda\} = \{\mathbf{x}: T(\mathbf{X}) \geq \tilde{\lambda}\}. \quad (14)$$

Заметим, что $\tilde{\lambda} = c$ подходит под условие. Следовательно, для любого $\theta_1 > \theta_0$

$$P_{\theta_1}(\mathbf{X} \in R) \leq P_{\theta_1}(\mathbf{X} \in S_{\lambda}) = P_{\theta_1}(T(\mathbf{X}) \geq c) = P_{\theta_1}(\mathbf{X} \in S). \quad (15)$$

Следовательно, S является равномерно наиболее мощным критерием. \square

Следствие. S является равномерно наиболее мощным критерием уровня доверия γ для проверки гипотезы $H_0: \theta = \theta_0$ против гипотезы $H_1: \theta > \theta_0$.

Теперь рассмотрим пример задачи на эту теорему.

Задача 25. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из распределения Бернулли $\text{Bern}(\theta)$, $\theta \in (0, 1)$. Найдите равномерно наиболее мощный критерий уровня значимости γ для проверки

- a) гипотеза $H_0: \theta \leq \theta_0$ против гипотезы $H_1: \theta > \theta_0$;
- b) гипотеза $H_0: \theta \geq \theta_0$ против гипотезы $H_1: \theta < \theta_0$.

Доказательство. Для начала запишем функцию правдоподобия для выборки:

$$p_{\theta}(\mathbf{X}) = \prod_{i=1}^n p_{\theta}(X_i) = \prod_{i=1}^n \theta^{X_i} (1 - \theta)^{1 - X_i} = \theta^{\sum_{i=1}^n X_i} (1 - \theta)^{n - \sum_{i=1}^n X_i}. \quad (16)$$

Вычислим отношение правдоподобия для $\theta_1 \neq \theta_0$:

$$\frac{p_{\theta_1}(\mathbf{X})}{p_{\theta_0}(\mathbf{X})} = \left(\frac{\theta_1}{\theta_0} \right)^{\sum_{i=1}^n X_i} \left(\frac{1 - \theta_1}{1 - \theta_0} \right)^{n - \sum_{i=1}^n X_i} \quad (17)$$

$$= \left(\frac{1 - \theta_1}{1 - \theta_0} \right)^n \left(\frac{\theta_1(1 - \theta_0)}{\theta_0(1 - \theta_1)} \right)^{\sum_{i=1}^n X_i}. \quad (18)$$

Теперь будем искать критерии.

- а) Если θ_1 из альтернативы, то $\theta_1 > \theta_0$ и $\theta_1(1 - \theta_0) > \theta_1(1 - \theta_1)$. Следовательно, при росте $\sum_{i=1}^n X_i$ отношение правдоподобия возрастает. Тогда искомая статистика $T(\mathbf{X}) = \sum_{i=1}^n X_i$ и критерий имеет вид $S = \{\mathbf{x} : \sum_{i=1}^n x_i \geq c\}$. Теперь нужно найти c :

$$P_{\theta_0}\left(\sum_{i=1}^n x_i \geq c\right) = \gamma \implies c = u_{1-\gamma}, \quad (19)$$

где $u_{1-\gamma}$ есть $(1 - \gamma)$ -квантиль биномиального распределения $\text{Bin}(n, \theta_0)$.

- б) В таком случае $\theta_1 < \theta_0$ и отношение правдоподобия уменьшается с ростом $\sum_{i=1}^n X_i$. Здесь ситуация ровно такая же, только нужно инвертировать знаки, то есть критерий имеет вид $S = \{\mathbf{x} : \sum_{i=1}^n x_i \leq c\}$ и $c = u_\gamma$.

□

12 Лекция 12

12.1 Критерий согласия

Рассмотрим так называемые критерии согласия. Такие критерии предлагают согласиться или же не согласиться с чем-то, что имеет простую природу.

Определение 47. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — функция распределения с неизвестной функцией распределения $F(x)$. Далее, введём гипотезы $H_0: F = F_0$ и $H_1: F \neq F_0$. Критерий для проверки гипотезы H_0 против альтернативы H_1 называется *критерием согласия*.

Что мы хотели бы видеть от такого критерия? В идеале увидеть что-то наподобие равномерно наиболее мощного критерия. Но эта задача почти всегда полностью безнадежна. Даже для хороших распределений, для которых выполнено свойство монотонного отношения правдоподобий, такого критерия может не быть. Например, для распределений $\text{Bin}(1, \theta)$, $\text{Exp}(\theta)$, $\mathcal{N}(\theta, 1)$ найти равномерно наиболее мощный критерий для проверки гипотезы $H_0: \theta = \theta_0$ против альтернативы $H_1: \theta \neq \theta_0$ не получится. С другой стороны, для некоторых распределений, для которых свойство монотонного отношения правдоподобий не выполнено, можно найти равномерно наиболее мощный критерий (например, $U[0, \theta]$). Тем самым такая задача в некоторой степени осмысленна и для некоторых, казалось бы, сложных распределений её можно решить.

Для того, чтобы её решить, пока что приходится всё делать по-честному. Поэтому перейдём задаче поиска *асимптотических критериев*. Такая постановка осмысленна: если не можем решить точно, то хотя бы решим приближенно. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка растущего размера. Мы хотели бы получить последовательность критериев $S_n = S_n(X_1, \dots, X_n)$ для проверки гипотезы $H_0: \theta = \theta_0$ против альтернативы $H_1: \theta \neq \theta_0$ такую, что

- Предел ошибки вероятности ошибки первого рода стремится к $\gamma \in (0, 1)$: $P_{\theta_0}(\mathbf{X} \in S_n) \rightarrow \gamma$.
- Последовательность критериев будет состоятельна: для любого $\theta \neq \theta_0$ $P_\theta(\mathbf{X} \in S_n) \rightarrow 1$.

Не смотря на то, что условия кажутся сложными для выполнения, оказывается, что решения есть. Рассмотрим один из таких критериев — критерий хи-квадрат Пирсона. Пусть есть выборка $\mathbf{X} = (X_1, \dots, X_n)$ из распределения P с конечным числом значений $\{a_1, \dots, a_m\}$:

$$p_j = P(X_1 = a_j) > 0, \quad j = 1, \dots, m, \quad \sum_{i=1}^m p_i = 1. \quad (20)$$

Введём вектор вероятностей $\mathbf{p} = (p_1, \dots, p_m)$. Гипотеза $H_0: F = F_0$ бужет сводиться к равенству векторов вероятностей, то есть задача согласия сводится к проверке гипотезы $H_0: \mathbf{p} = \mathbf{p}^0$ против альтернативы $H_1: \mathbf{p} \neq \mathbf{p}^0$. Для проверки данной гипотезы вводится *статистика хи-квадрат*:

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{j=1}^m \frac{(\mu_j - np_j^0)^2}{np_j^0}, \quad \text{где } \mu_j = \sum_{i=1}^n \mathbf{1}_{X_i = a_j}. \quad (21)$$

Критерий же устроен следующим образом: пусть $u_{1-\gamma} - (1 - \gamma)$ -квантиль распределения χ_{m-1}^2 . Тогда гипотеза H_0 отвергается тогда и только тогда, когда $\hat{\chi}_n^2(\mathbf{X}) > u_{1-\gamma}$. Теперь докажем, что это хороший асимптотический критерий.

Утверждение. Критерий хи-квадрат состоятелен.

Доказательство. Для начала немного перепишем статистику хи-квадрат:

$$\hat{\chi}_n^2(\mathbf{X}) = n \sum_{j=1}^m \frac{1}{p_j^0} \left(\frac{\mu_j}{n} - p_j^0 \right)^2. \quad (22)$$

Далее, если $\mathbf{p} \neq \mathbf{p}^0$, то найдётся индекс $i \in \{1, 2, \dots, m\}$ такой, что $p_j \neq p_j^0$. Однако по усиленному закону больших чисел $\mu_j/n \rightarrow p_j$ почти наверное. Отсюда следует, что $\hat{\chi}_n^2(\mathbf{X})$ линейно стремится к бесконечности, так как

$$\frac{\hat{\chi}_n^2(\mathbf{X})}{n} \xrightarrow{\text{п.н.}} \sum_{j=1}^m \frac{(p_j - p_j^0)^2}{p_j^0} > 0. \quad (23)$$

Тем самым критерий состоятелен: $P(\hat{\chi}_n^2(\mathbf{X}) > u_{1-\gamma}) \rightarrow 1$. \square

Первое свойство будет получаться из следующей теоремы:

Теорема 27 (Пирсона). *Если выполнена гипотеза H_0 , то имеет место следующая сходимость по распределению:*

$$\hat{\chi}_n^2(\mathbf{X}) \xrightarrow[n \rightarrow \infty]{d} \chi_{m-1}^2. \quad (24)$$

Доказательство. Пусть $\mathbf{p} = \mathbf{p}^0$. Для всех $i = 1, \dots, n$ введём случайный вектор $\mathbf{Y}_i = (\mathbf{1}_{X_i=a_1}, \dots, \mathbf{1}_{X_i=a_m})$. Понятно, что такие векторы независимы и одинаково распределены, причём $E[\mathbf{Y}_i] = \mathbf{p}^0$. Далее заметим, что

$$\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} = \frac{1}{n} \begin{pmatrix} \mu_1 \\ \dots \\ \mu_m \end{pmatrix}. \quad (25)$$

Следовательно, по многомерной центральной предельной теореме

$$\sqrt{n} \left(\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} - \mathbf{p}^0 \right) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}), \quad (26)$$

где $\mathbf{\Sigma} = D[\mathbf{Y}_i]$. Заметим, что

$$\text{cov}(\mathbf{1}_{X_1=a_i}, \mathbf{1}_{X_1=a_j}) = P(X_1 = a_i, X_1 = a_j) - p_i^0 p_j^0 = \begin{cases} p_i^0 - p_i^0 p_j^0, & i = j \\ -p_i^0 p_j^0, & i \neq j \end{cases} \quad (27)$$

Тогда запишем матрицу ковариаций следующим образом: $\mathbf{\Sigma} = \mathbf{B} - \mathbf{p}^0(\mathbf{p}^0)^\top$, где $\mathbf{B} = \text{diag}(\mathbf{p}^0)$. Теперь наша цель — поворотать вектор в ЦПТ так, чтобы матрица ковариаций стала разумной: нам нужно получить распределение хи-квадрат, поэтому желательно получить матрицу, чем-то похожую на матрицу стандартного нормального распределения). Пусть

$$\boldsymbol{\xi}'_n = \sqrt{n} \left(\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} - \mathbf{p}^0 \right). \quad (28)$$

Тогда по центральной предельной теореме

$$(\sqrt{\mathbf{B}})^{-1} \boldsymbol{\xi}'_n \xrightarrow[n \rightarrow \infty]{d} (\sqrt{\mathbf{B}})^{-1} \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}) = \mathcal{N}(\mathbf{0}, (\sqrt{\mathbf{B}})^{-1} \mathbf{\Sigma} (\sqrt{\mathbf{B}})^{-1}). \quad (29)$$

Новая матрица ковариаций равна $\mathbf{I}_m - \mathbf{z}\mathbf{z}^\top$, где $\mathbf{z} = (\sqrt{\mathbf{B}})^{-1} \mathbf{p}^0 = (\sqrt{p_1^0}, \dots, \sqrt{p_m^0})$. Заметим, что вектор \mathbf{z} единичный, а это означает, что его можно перевести поворотами в базисный вектор. Рассмотрим ортогональную матрицу \mathbf{C} , у которой первая строка равна \mathbf{z} , а остальные равны чему угодно. Тогда по теореме о наследовании сходимости

$$\boldsymbol{\xi}''_n = \mathbf{C}(\sqrt{\mathbf{B}})^{-1} \boldsymbol{\xi}'_n \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(\mathbf{0}, \mathbf{C}(\mathbf{I}_m - \mathbf{z}\mathbf{z}^\top) \mathbf{C}^\top). \quad (30)$$

Однако $\mathbf{C}\mathbf{z} = (1, 0, \dots, 0)^\top$, а $\mathbf{C}\mathbf{C}^\top = \mathbf{I}_m$. Поэтому $\mathbf{C}(\mathbf{I}_m - \mathbf{z}\mathbf{z}^\top) \mathbf{C}^\top = \text{diag}(0, 1, \dots, 1) = \mathbf{I}'_m$. Снова по теореме о наследовании сходимости

$$\|\boldsymbol{\xi}''_n\|^2 \xrightarrow[n \rightarrow \infty]{d} \|\mathcal{N}(\mathbf{0}, \mathbf{I}'_m)\|^2 \sim \chi_{m-1}^2. \quad (31)$$

Но, как известно, помножение вектора на ортогональную матрицу не изменяет его норму. Тогда

$$\|(\sqrt{\mathbf{B}})^{-1}\boldsymbol{\xi}'_n\|^2 \xrightarrow[n \rightarrow \infty]{d} \chi_{m-1}^2. \quad (32)$$

Осталось заметить, что

$$(\sqrt{\mathbf{B}})^{-1}\boldsymbol{\xi}'_n = (\sqrt{\mathbf{B}})^{-1}\sqrt{n} \begin{pmatrix} (\mu_1 - np_1^0)/n \\ \dots \\ (\mu_m - np_m^0)/n \end{pmatrix} = \begin{pmatrix} (\mu_1 - np_1^0)/\sqrt{np_1^0} \\ \dots \\ (\mu_m - np_m^0)/\sqrt{np_m^0} \end{pmatrix} \quad (33)$$

Тем самым $\hat{\chi}_n^2(\mathbf{X}) = \|(\sqrt{\mathbf{B}})^{-1}\boldsymbol{\xi}'_n\|^2$. \square

Насколько этот метод применим? На практике считается, что должно быть следующее ограничение: $np_j^0 \geq 5$ для всех $j = 1, \dots, m$.

12.2 Параметрический хи-квадрат

Метод достаточно мощный, хоть мы и пожертвовали тем, что вероятность ошибки первого рода не равна γ , а только стремится к ней. Что самое ироничное: отказ от точности дал некоторые хорошие свойства. Например, его можно применять, когда модель параметрическая. В таком случае критерий называется *параметрическим критерием хи-квадрат*.

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ — выборка из неизвестного распределения со значениями $\{a_1, \dots, a_m\}$:

$$p_j(\theta) = P_\theta(X_1 = a_j), \quad \theta \in \Theta \quad (34)$$

В данном случае гипотеза и альтернатива формулируются следующим образом: $H_0: \mathbf{p} \in \{P_\theta, \theta \in \Theta\}$ против $H_1: \mathbf{p} \notin \{P_\theta, \theta \in \Theta\}$. Теперь хотелось бы составить статистику хи-квадрат, но особо не получится, так как \mathbf{p}^0 больше нет. Можно ли его чем-то заменить? Заменим на оценку:

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{j=1}^m \frac{(\mu_j - n\hat{p}_j(\mathbf{X}))^2}{n\hat{p}_j(\mathbf{X})}. \quad (35)$$

Какую оценку взять? Самым разумным выбором является оценка максимального правдоподобия: $\hat{p}_j(\mathbf{X}) = p_j(\hat{\theta}(\mathbf{X}))$, где $\hat{\theta}(\mathbf{X})$ — оценка максимального правдоподобия для θ . Оказывается, что в таком случае можно доказать результат, похожий на теорему Пирсона: если $\theta \in \mathbb{R}^k$, $k < m - 1$, то при выполнении определённых условий

$$\hat{\chi}_n^2(\mathbf{X}) \xrightarrow[n \rightarrow \infty]{d} \chi_{m-1-k}^2. \quad (36)$$

Данную теорему мы докажем на следующей лекции.

13 Лекция 13

Вернёмся к параметрическому критерию хи-квадрат. Ранее мы сказали, что при выполнении определённых условий будет выполнен аналог теоремы Пирсона. Сформулируем их:

- $\Theta \subseteq \mathbb{R}^s$, $s < m$ — открытое множество.
- Для любого значения параметра все вероятности отделены от нуля: $p_i(\theta) \geq c^2 > 0$.
- Будем считать, что $\frac{\partial p_i(\theta)}{\partial \theta_j}$ и $\frac{\partial^2 p_i(\theta)}{\partial \theta_j \partial \theta_k}$ непрерывны на всём Θ .
- Матрица $\mathbf{D} = \left\| \frac{\partial p_i(\theta)}{\partial \theta_j} \right\|_{i,j=1}^{m,s}$ имеет ранг s для любого $\theta \in \Theta$.

Теперь можно сформулировать теорему про параметрический критерий хи-квадрат.

Теорема 28. Пусть выполнены условия модели. Введём следующую систему уравнений:

$$\sum_{i=1}^m \frac{\mu_i}{p_i(\theta)} \frac{\partial p_i(\theta)}{\partial \theta_j} = 0, \quad j = 1, \dots, s, \quad \mu_j = \sum_{i=1}^n \mathbf{1}_{X_i=a_j}. \quad (37)$$

Если верна гипотеза H_0 , то с вероятностью, стремящейся к 1, данная система имеет единственное решение $\hat{\theta}(\mathbf{X})$ такое, что $\hat{\theta}(\mathbf{X})$ сходится по вероятности к истинному значению параметра θ и

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{j=1}^m \frac{(\mu_j - np_j(\hat{\theta}(\mathbf{X})))^2}{np_j(\hat{\theta}(\mathbf{X}))} \xrightarrow[n \rightarrow \infty]{d} \chi_{m-1-s}^2. \quad (38)$$

Доказательство. Пусть θ_0 — истинное значение параметра θ . Далее, для удобства скажем, что

$$p_i(\theta_0) = p_i^0, \quad \left. \frac{\partial p_i(\theta)}{\partial \theta_j} \right|_{\theta=\theta_0} = \left(\frac{\partial p_i}{\partial \theta_j} \right)_0.$$

Теперь введём матрицы

$$\mathbf{P}_0 = \text{diag} \left(\frac{1}{\sqrt{p_1^0}}, \dots, \frac{1}{\sqrt{p_m^0}} \right) \in \mathbb{R}^{m \times m}, \quad \mathbf{D}_0 = \mathbf{D}(\theta_0) = \left\| \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 \right\|_{i,j=1}^{m,s} \in \mathbb{R}^{m \times s}, \quad \mathbf{B}_0 = \mathbf{P}_0 \mathbf{D}_0.$$

Теперь выразим систему из условия через данные матрицы. Введём ещё одно обозначение:

$$\begin{aligned} \omega_j(\theta) = & \sum_{i=1}^m \frac{\mu_i - np_i^0}{n} \left(\frac{1}{p_i} \frac{\partial p_i(\theta)}{\partial \theta_j} - \frac{1}{p_i^0} \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 \right) - \sum_{i=1}^m (p_i - p_i^0) \left(\frac{1}{p_i} \frac{\partial p_i(\theta)}{\partial \theta_j} - \frac{1}{p_i^0} \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 \right) \\ & - \sum_{i=1}^m \frac{1}{p_i^0} \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 \left(p_i - p_i^0 - \sum_{k=1}^s \left(\frac{\partial p_i}{\partial \theta_k} \right)_0 (\theta_k - \theta_k^0) \right). \end{aligned}$$

Зачем вводить это безобразие? Для того, чтобы доказать, что оно мало. Воспользуемся системой из условия и тем, что сумма p_i равна единице:

$$\omega_j(\theta) = - \sum_{i=1}^m \frac{\mu_i - np_i^0}{n} \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 + \sum_{i=1}^m \sum_{k=1}^s \frac{1}{p_i^0} \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 \left(\frac{\partial p_i}{\partial \theta_k} \right)_0 (\theta_k - \theta_k^0).$$

Несложно показать, что если выполнено такое равенство, то верна система из условия. Теперь перепишем это равенство в матричном виде. Введём вектор \mathbf{Y}_j по следующему правилу: $Y_j = (\mu_j - np_j^0) / \sqrt{np_j^0}$. Тогда (проверьте!)

$$\omega(\theta) = - \frac{1}{\sqrt{n}} \mathbf{B}_0^\top \mathbf{Y} + \mathbf{B}_0^\top \mathbf{B}_0 (\theta - \theta_0).$$

Далее, заметим, что матрица $\mathbf{B}_0^\top \mathbf{B}_0$ обратима, так как \mathbf{D}_0 имеет ранг s . В таком случае

$$\theta - \theta_0 = \frac{1}{\sqrt{n}} (\mathbf{B}_0^\top \mathbf{B}_0)^{-1} \mathbf{B}_0^\top \mathbf{Y} + (\mathbf{B}_0^\top \mathbf{B}_0)^{-1} \omega(\theta).$$

Теперь нужно доказать несколько фактов.

Лемма. Пусть $\lambda = \lambda(n) \rightarrow \infty$, $\lambda^2 = o(n)$. Тогда с вероятностью $1 - \lambda^{-2}$ для всех $i \in \{1, \dots, m\}$

$$|Y_i| \leq \frac{\lambda}{c}.$$

Доказательство. Заметим, что $\mu_i \sim \text{Bin}(n, p_i^0)$. Тогда по неравенству Чебышева

$$\mathbf{P}_{\theta_0} (|\mu_i - np_i^0| \geq \lambda \sqrt{n}) \leq \frac{\mathbf{D}_{\theta_0}[\mu_i]}{n \lambda^2} \leq \frac{p_i^0}{\lambda^2}.$$

Следовательно,

$$\mathbf{P}_{\theta_0} (\exists i: |\mu_i - np_i^0| \geq \lambda \sqrt{n}) \leq \sum_{i=1}^m \mathbf{P}_{\theta_0} (|\mu_i - np_i^0| \geq \lambda \sqrt{n}) \leq \sum_{i=1}^m \frac{p_i^0}{\lambda^2} = \frac{1}{\lambda^2}.$$

Тогда, с вероятностью не меньше $1 - \lambda^{-2}$ для всех i

$$|Y_i| = \frac{|\mu_i - np_i^0|}{\sqrt{np_i^0}} \leq \frac{\lambda \sqrt{n}}{\sqrt{np_i^0}} \leq \frac{\lambda}{c}. \quad \square$$

Лемма.

$$|\omega_j(\theta') - \omega_j(\theta'')| \leq K_1 \|\theta' - \theta''\| \left(\|\theta' - \theta_0\| + \|\theta'' - \theta_0\| + \frac{\lambda}{\sqrt{n}} \right),$$

где K_1 — некоторая абсолютная константа.

Доказательство. Оставляется читателю в качестве упражнения.¹⁸ *Указание:* Рассмотрим разность (сразу учтём, что p_i суммируются в 1):

$$\begin{aligned} w_j(\boldsymbol{\theta}') - w_j(\boldsymbol{\theta}'') &= \sum_{i=1}^m \frac{\mu_i - np_i^0}{n} \left(\frac{1}{p_i(\boldsymbol{\theta}')} \frac{\partial p_i(\boldsymbol{\theta}')}{\partial \theta'_j} - \frac{1}{p_i(\boldsymbol{\theta}'')} \frac{\partial p_i(\boldsymbol{\theta}'')}{\partial \theta''_j} \right) \\ &+ \sum_{i=1}^m \frac{p_i(\boldsymbol{\theta}') - p_i(\boldsymbol{\theta}'')}{p_i^0} \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 + \sum_{i=1}^m p_i^0 \left(\frac{1}{p_i(\boldsymbol{\theta}')} \frac{\partial p_i(\boldsymbol{\theta}')}{\partial \theta'_j} - \frac{1}{p_i(\boldsymbol{\theta}'')} \frac{\partial p_i(\boldsymbol{\theta}'')}{\partial \theta''_j} \right) \\ &- \sum_{i=1}^m \frac{1}{p_i^0} \left(\frac{\partial p_i}{\partial \theta_j} \right)_0 \left(p_i(\boldsymbol{\theta}') - p_i(\boldsymbol{\theta}'') - \sum_{k=1}^s \left(\frac{\partial p_i}{\partial \theta_k} \right)_0 (\theta'_k - \theta''_k) \right). \end{aligned}$$

Далее, воспользуемся разложением $p_i(\boldsymbol{\theta}')$ и $p_i(\boldsymbol{\theta}'')$ в ряд Тейлора до второго порядка в точке $\boldsymbol{\theta}_0$. \square

Теперь осталось построить саму оценку. Рассмотрим уравнение

$$\boldsymbol{\theta} = \boldsymbol{\theta}_0 + \frac{1}{\sqrt{n}} (\mathbf{B}_0^\top \mathbf{B}_0)^{-1} \mathbf{B}_0^\top \mathbf{Y} + (\mathbf{B}_0^\top \mathbf{B}_0)^{-1} \boldsymbol{\omega}(\boldsymbol{\theta}).$$

Заметим, что $\boldsymbol{\omega}(\boldsymbol{\theta}_0) = 0$. Далее, определим последовательность приближений $(\hat{\boldsymbol{\theta}}_l, l \in \mathbb{N})$ по следующему правилу:

$$\begin{aligned} \hat{\boldsymbol{\theta}}_1 &= \boldsymbol{\theta}_0 + \frac{1}{\sqrt{n}} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{Y}, \\ \hat{\boldsymbol{\theta}}_l &= \hat{\boldsymbol{\theta}}_1 + (\mathbf{B}^\top \mathbf{B})^{-1} \boldsymbol{\omega}(\hat{\boldsymbol{\theta}}_{l-1}). \end{aligned}$$

Заметим, что $\|\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_0\| \leq K_2 \lambda / \sqrt{n}$ для какой-то константы K_2 . \square

14 Лекция 14

15 Лекция 15

15.1 Слабая сходимость

Определение 48. Метрическим пространством будем называть (X, d) , где X есть некоторое множество, а $d: X \times X \mapsto \mathbb{R}_+$ — функция, называемая метрикой, и обладающая следующими свойствами:

- a) Аксиома тождества: $d(x, x) = 0$ тогда и только тогда, когда $x = 0$,
- b) Аксиома симметрии: $d(x, y) = d(y, x)$,
- c) Неравенство треугольника: $d(x, z) \leq d(x, y) + d(y, z)$.

Определение 49. Пусть (S, ρ) — метрическое пространство. Минимальную сигма-алгебру, порождённую открытыми шарами в S , будем называть борелевской сигма-алгеброй $\mathcal{B}(S)$.

Определение 50. Пусть задано метрическое пространство (S, ρ) и последовательность $\{Q_n\}_{n=1}^\infty$ вероятностных мер на S . Будем говорить, что Q_n слабо сходятся к вероятностной мере Q на S , если для любой ограниченной непрерывной функции $f: S \mapsto \mathbb{R}$

$$\lim_{n \rightarrow \infty} \int_S f(x) Q_n(dx) = \int_S f(x) Q(dx). \quad (39)$$

Обозначение: $Q_n \xrightarrow{W} Q$.

Теорема 29 (Александров). Пусть $\{Q_n\}_{n=1}^\infty$ и Q — вероятностные меры на метрическом пространстве (S, ρ) . Тогда следующие утверждения эквивалентны:

- a) $Q_n \xrightarrow{W} Q$,
- b) Для любого замкнутого множества $F \subset S$ $\lim_{n \rightarrow \infty} Q_n(F) \leq Q(F)$,
- c) Для любого открытого множества $G \subset S$ $\lim_{n \rightarrow \infty} Q_n(G) \geq Q(G)$,

¹⁸Если кто сможет довести эту бандуру до конца — пинганите в Telegram.

d) Для любого борелевского множества $B \in \mathcal{B}(S)$ такого, что $Q(\partial B) = 0$, $Q_n(B) \rightarrow Q(B)$ при $n \rightarrow \infty$.

Доказательство. Для начала докажем, что если выполнена слабая сходимость, то для произвольного замкнутого множества F выполнено ограничение снизу. Для произвольного $\varepsilon > 0$ введём функцию

$$f_\varepsilon(x) = \left(1 - \frac{\rho(x, F)}{\varepsilon}\right)^+, \quad \text{где } \rho(x, F) = \inf_{y \in F} d(x, y). \quad (40)$$

Заметим, что она непрерывна и ограничена. Тогда верна следующая цепочка:

$$\overline{\lim}_{n \rightarrow \infty} Q_n(F) = \overline{\lim}_{n \rightarrow \infty} \int_S \mathbf{1}_{x \in F} Q_n(dx) \leq \overline{\lim}_{n \rightarrow \infty} \int_S f_\varepsilon(x) Q_n(dx) = \int_S f_\varepsilon(x) Q(dx) \leq Q(F^\varepsilon), \quad (41)$$

где $F^\varepsilon = \{x \in S: \rho(x, F) \leq \varepsilon\}$. Осталось заметить, что F^ε монотонно сжимается к F при $\varepsilon \rightarrow 0$. Тогда по непрерывности вероятностной меры $Q(F^\varepsilon) \rightarrow Q(F)$. Следовательно, при устремлении ε к нулю получаем второй пункт.

Далее, эквивалентность второго и третьего пунктов очевидна, так как от одного к другому можно переходить, рассматривая дополнения множеств. Теперь докажем, что из второго и третьих пунктов следует четвёртый пункт. Для этого возьмём произвольное борелевское множество $B \in \mathcal{B}(S)$ такое, что $Q(\partial B) = 0$. Теперь введём два множества: $F = [B]$ и $G = B \setminus \partial B$. Заметим, что F замкнуто, а G открыто, причём их меры равны мере B (так как мера границы есть ноль). Тогда несложно показать, что искомый предел существует и равен $Q(B)$:

$$\overline{\lim}_{n \rightarrow \infty} Q_n(B) \leq \overline{\lim}_{n \rightarrow \infty} Q_n(F) \leq Q(F) = Q(B), \quad (42)$$

$$\underline{\lim}_{n \rightarrow \infty} Q_n(B) \geq \underline{\lim}_{n \rightarrow \infty} Q_n(G) \geq Q(G) = Q(B). \quad (43)$$

Осталось доказать, что из последнего пункта следует первый. Возьмём некоторую ограниченную непрерывную функцию $f: S \rightarrow \mathbb{R}$ такую, что $|f(x)| < M$ для всех $x \in S$. Далее, построим следующее множество:

$$D = \{t \in [-M, M]: Q(\{x \in S: f(x) = t\}) > 0\}. \quad (44)$$

Данное множество не более, чем счётно. Теперь зафиксируем произвольное натуральное k и возьмём разбиение $-M = t_0 < t_1 < \dots < t_k = M$ отрезка $[-M, M]$ такое, что ни одно из t_i не содержится в D . Далее, построим набор множеств $\{B_i\}_{i=1}^k$ по следующему правилу: $B_i = \{x \in S: t_{i-1} \leq f(x) < t_i\}$. Заметим, что $\partial B_i \subseteq f^{-1}(\{t_{i-1}\}) \cup f^{-1}(\{t_i\})$. Но оба прообраза имеют нулевую меру, поэтому граница B_i тоже имеет нулевую меру. Следовательно, $Q_n(B_i) \rightarrow Q(B_i)$ при $n \rightarrow \infty$ для всех $i = 1, \dots, k$. Теперь рассмотрим следующий верхний предел:

$$\Delta = \overline{\lim}_{n \rightarrow \infty} \left| \int_S f(x) Q_n(dx) - \int_S f(x) Q(dx) \right|. \quad (45)$$

Ограничим её сверху суммой трёх верхних пределов:

$$\Delta \leq \overline{\lim}_{n \rightarrow \infty} \left| \int_S f(x) Q_n(dx) - \sum_{i=1}^k t_{i-1} Q_n(B_i) \right| \quad (46)$$

$$+ \overline{\lim}_{n \rightarrow \infty} \left| \sum_{i=1}^k t_{i-1} Q_n(B_i) - \sum_{i=1}^k t_{i-1} Q(B_i) \right| \quad (47)$$

$$+ \overline{\lim}_{n \rightarrow \infty} \left| \int_S f(x) Q(dx) - \sum_{i=1}^k t_{i-1} Q(B_i) \right| \quad (48)$$

Второй предел сразу равен нулю. Теперь покажем, что первый (а так же и третий) предел не превосходит $\max_{i=1, \dots, k} |t_i - t_{i-1}|$. Для этого заметим, что

$$\int_S f(x) Q_n(dx) = \sum_{i=1}^k \int_{B_i} f(x) Q_n(dx) \leq \sum_{i=1}^k t_i Q_n(B_i). \quad (49)$$

Следовательно,

$$\left| \int_S f(x) Q_n(dx) - \sum_{i=1}^k t_{i-1} Q_n(B_i) \right| \leq \left| \sum_{i=1}^k (t_i - t_{i-1}) Q_n(B_i) \right| \leq \max_{i=1, \dots, k} |t_i - t_{i-1}|. \quad (50)$$

Отсюда получаем, что

$$\Delta \leq 2 \max_{i=1, \dots, k} |t_i - t_{i-1}|. \quad (51)$$

Но этот максимум стремится к нулю при устремлении диаметра разбиения к нулю. Следовательно, имеет место слабая сходимость. \square

15.2 Случайные процессы

Пусть $X = (X_t, t \in [0, 1])$ — случайный процесс. Тогда X — случайный элемент со значением в пространстве функций на $[0, 1]$ с цилиндрической сигма-алгеброй.

Определение 51. Пространство функций на $[0, 1]$ — это множество $\mathbb{R}^{[0,1]} = \{y = (y(t), t \in [0, 1]), y(t) \in \mathbb{R}\}$.

Определение 52. Пусть $B \in \mathcal{B}(\mathbb{R})$ и $t \in [0, 1]$. *Элементарным цилиндром* называется множество $C(t, B) = \{y \in \mathbb{R}^{[0,1]} : y(t) \in B\}$.

Определение 53. *Цилиндрической сигма-алгеброй* называется минимальная сигма-алгебра, порождённая элементарными цилиндрами:

$$\text{Cyl}(\mathbb{R}^{[0,1]}) = \sigma(\{C(t, B) : t \in [0, 1], B \in \mathcal{B}(\mathbb{R})\}) \quad (52)$$

Утверждение. Отображение $X : \Omega \mapsto (\mathbb{R}^{[0,1]}, \text{Cyl}(\mathbb{R}^{[0,1]}))$ является измеримым.

Из этого следует, что имеет смысл вводить распределение случайного процесса $P_X(C) = P(X \in C)$, $C \in \text{Cyl}(\mathbb{R}^{[0,1]})$. Но пространство $\mathbb{R}^{[0,1]}$ плохое. Поэтому обычно его сужают до пространства непрерывных функций из $[0, 1]$ в \mathbb{R} , которое обозначают через $\mathcal{C}[0, 1]$. На нём есть норма $\|x\| = \max_{t \in [0,1]} |x(t)|$. Следовательно, можно построить борелевскую сигма-алгебру $\mathcal{B}(\mathcal{C}[0, 1])$.

Лемма. Для пространства $\mathcal{C}[0, 1]$ борелевская сигма-алгебра совпадает с цилиндрической сигма-алгеброй: $\mathcal{B}(\mathcal{C}[0, 1]) = \text{Cyl}(\mathcal{C}[0, 1])$.

Доказательство. Для начала докажем, что выполнено следующее вложение: $\mathcal{B}(\mathcal{C}[0, 1]) \subseteq \text{Cyl}(\mathcal{C}[0, 1])$. Для этого возьмём замкнутый шар $B_r[x] = \{y \in \mathcal{C}[0, 1] : \|y - x\| \leq r\}$. Тогда по непрерывности можно сказать, что

$$B_r[x] = \bigcap_{t \in \mathbb{Q}} \{y \in \mathcal{C}[0, 1] : |y(t) - x(t)| \leq r\}. \quad (53)$$

Но это есть ни что иное, как счётное пересечение элементарных цилиндров. Следовательно, $B_r[x] \in \text{Cyl}(\mathcal{C}[0, 1])$.

Теперь докажем обратное вложение. Для этого рассмотрим следующий цилиндр:

$$C(t, [a, b]) = \{y \in \mathcal{C}[0, 1] : y(t) \in [a, b]\}. \quad (54)$$

Теперь введём функцию $h_t : \mathcal{C}[0, 1] \mapsto \mathbb{R}$, действующую по правилу $h_t(x) = x(t)$. Несложно понять, что это непрерывная функция. Теперь заметим, что $C(t, [a, b]) = h_t^{-1}([a, b])$. Но это будет замкнутое множество, так как для непрерывной функции прообраз замкнутого множества замкнут. Следовательно, $C(t, [a, b]) \in \mathcal{B}(\mathcal{C}[0, 1])$. \square

Определение 54. Пусть $X^{(n)} = (X_t^{(n)}, t \in [0, 1])$ — последовательность случайных процессов. Будем говорить, что она сходится по распределению к случайному процессу $X = (X_t, t \in [0, 1])$, если имеет место слабая сходимост распределений: $P_{X^{(n)}} \xrightarrow{W} P_X$ при $n \rightarrow \infty$. Обозначение: $X^{(n)} \xrightarrow{D} X$.

Теорема 30 (о наследовании сходимости). Пусть $h : \mathcal{C}[0, 1] \mapsto \mathbb{R}$ — непрерывная функция и $X^{(n)} \xrightarrow{D} X$ при $n \rightarrow \infty$. Тогда $h(X^{(n)}) \xrightarrow{d} h(X)$.

Доказательство. Пусть $y : \mathbb{R} \mapsto \mathbb{R}$ — ограниченная непрерывная функция. Тогда композиция $g = y \circ h \mapsto \mathcal{C}[0, 1] \mapsto \mathbb{R}$ будет ограниченной непрерывной функцией. Рассмотрим $E[y(h(X^{(n)}))]$. Оно равно

$$E[y(h(X^{(n)}))] = E[g(X^{(n)})] = \int_{\mathcal{C}[0,1]} g(x) P_{X^{(n)}}(dx). \quad (55)$$

Воспользуемся слабой сходимостью распределений:

$$\lim_{n \rightarrow \infty} \int_{\mathcal{C}[0,1]} g(x) P_{X^{(n)}}(dx) = \int_{\mathcal{C}[0,1]} g(x) P_X(dx) = E[g(X)]. \quad (56)$$

Следовательно, для любой ограниченной непрерывной функции $y : \mathbb{R} \mapsto \mathbb{R}$ выполнена следующая сходимост: $E[y(h(X^{(n)}))] \rightarrow E[y(h(X))]$ при $n \rightarrow \infty$. Но это означает, что $h(X^{(n)}) \xrightarrow{d} h(X)$. \square

15.3 Назад к критерию Колмогорова

Теорема 31 (принцип инвариантности Донскера-Прохорова). Пусть $\{\xi_n\}_{n=1}^\infty$ — последовательность независимых и одинаково распределённых случайных величин таких, что $E[\xi_1] = 0$ и $D[\xi_1] = 1$. Введём последовательность случайных величин $\{S_n\}_{n=0}^\infty$ по следующему правилу: $S_0 = 0$, $S_n = \xi_1 + \dots + \xi_n$. Далее, построим случайный процесс $X^{(n)} = (X_t^{(n)}, t \in [0, 1])$, как линейную интерполяцию S_0, \dots, S_n :

$$X_t^{(n)} = \frac{S_k}{\sqrt{n}}(k+1-nt) + \frac{S_{k+1}}{\sqrt{n}}(nt-k) \text{ при } t \in \left[\frac{k}{n}, \frac{k+1}{n}\right], \quad k = 0, \dots, n-1. \quad (57)$$

Тогда $X^{(n)} \xrightarrow{D} W$, где $W = (W_t, t \in [0, 1])$ — винеровский процесс.

Применим её к доказательству критерия Колмогорова.

Теорема 32. Пусть X_1, \dots, X_n — независимые и одинаково распределённые случайные величины с распределением $U[0, 1]$, а $x \in [0, 1]$. Далее, вводится статистика

$$D_n = \sup_{x \in [0, 1]} |\hat{F}_n(x) - x|. \quad (58)$$

Тогда имеет место следующая сходимость:

$$\sqrt{n}D_n \xrightarrow{d} \sup_{t \in [0, 1]} |W_t - tW_1|. \quad (59)$$

Доказательство. Для начала видоизменим определение D_n . Для этого поймём, где может достигаться супремум. Так как \hat{F}_n является кусочно-постоянной функцией, то супремум может достигаться только в точках разрыва, ибо иначе значение можно увеличить, немного уменьшив x — от этого значение $\hat{F}_n(x)$ не изменится. Следовательно,

$$D_n = \max_{k=1, \dots, n} \left| X_{(k)} - \frac{k}{n} \right|. \quad (60)$$

Ранее доказывалось, что если ξ_1, \dots, ξ_{n+1} — независимые и одинаково распределённые случайные величины с стандартным экспоненциальным распределением $\text{Exp}(1)$, а $S_k = \xi_1 + \dots + \xi_k$, то

$$(X_{(1)}, \dots, X_{(n)}) \stackrel{d}{=} \left(\frac{S_1}{S_{n+1}}, \dots, \frac{S_n}{S_{n+1}} \right). \quad (61)$$

Тогда распределение $\sqrt{n}D_n$ равно распределению

$$\sqrt{n} \max_{k=1, \dots, n} \left| \frac{S_k}{S_{n+1}} - \frac{k}{n} \right|. \quad (62)$$

Теперь заметим, что в пределе распределение выражения выше будет распределению

$$\Delta = \sqrt{n+1} \max_{k=1, \dots, n} \left| \frac{S_k}{S_{n+1}} - \frac{k}{n+1} \right|. \quad (63)$$

Преобразуем выражение:

$$\Delta = \frac{n+1}{S_{n+1}} \max_{k=1, \dots, n} \left| \frac{S_k}{\sqrt{n+1}} - \frac{kS_{n+1}}{(n+1)\sqrt{n+1}} \right| \quad (64)$$

$$= \frac{n+1}{S_{n+1}} \max_{k=1, \dots, n} \left| \frac{S_k - k}{\sqrt{n+1}} - \frac{k}{n+1} \frac{S_{n+1} - (n+1)}{\sqrt{n+1}} \right| = T_n. \quad (65)$$

Теперь введём случайный процесс $X^{(n)} = (X_t^{(n)}, t \in [0, 1])$ по следующему правилу:

$$X_t^{(n)} = \frac{S_k}{\sqrt{n+1}}(k+1-(n+1)t) + \frac{S_{k+1}}{\sqrt{n+1}}((n+1)t-k) \text{ при } t \in \left[\frac{k}{n+1}, \frac{k+1}{n+1}\right], \quad k = 0, \dots, n. \quad (66)$$

Далее заметим, что траектории случайного процесса $X^{(n)}$ — это ломаные. Поэтому

$$T_n \stackrel{d}{=} \frac{n+1}{S_{n+1}} \sup_{t \in [0, 1]} |X_t^{(n)} - tX_1^{(n)}|. \quad (67)$$

Согласно лемме Слуцкого, усиленному закону больших чисел, принципу инвариантности и теореме о наследовании сходимости получаем, что

$$\Delta \xrightarrow{d} \sup_{t \in [0, 1]} |W_t - tW_1|. \quad (68)$$

Тем самым получаем желаемое. \square

Определение 55. Случайный процесс $W^0 = (W_t^0, t \in [0, 1])$, где $W_t^0 = W_t - tW_1$, называется *броуновским мостом*.

Остался последний вопрос: как показать, что распределение супремума броуновского моста будет равно распределению Колмогорова? Это уже не самая тривиальная задача.

16 Лекция 16

17 Лекция 17

17.1 Ранговые методы

Вернёмся к критериям однородности, но несколько изменим формулировку. Пусть $\mathbf{X} = (X_1, \dots, X_m)$ и $\mathbf{Y} = (Y_1, \dots, Y_n)$ — независимые выборки с некоторыми неизвестными непрерывными распределениями. Можно ли сказать, что X_i равно Y_j по распределению для всех $i = 1, \dots, m, j = 1, \dots, n$?

Определение 56. Соберём выборки \mathbf{X} и \mathbf{Y} в вектор \mathbf{Z} . *Рангом* случайной величины Y_j называется

$$R(Y_j) = \sum_{i=1}^{n+m} \mathbf{1}_{Z_i \leq Y_j}. \quad (69)$$

Определение 57. *Статистикой ранговых сумм Вилкоксона* называется

$$W_{m,n}(\mathbf{X}, \mathbf{Y}) = \sum_{j=1}^n R(Y_j). \quad (70)$$

Утверждение. Если верна гипотеза однородности, то статистика ранговых сумм Вилкоксона не зависит от распределения выборок.

Доказательство. Пусть $F(x)$ — функция распределения элементов выборок. Так как она непрерывна, то $\mathbf{1}_{X_i \leq Y_j} = \mathbf{1}_{F(X_i) \leq F(Y_j)}$. Но $F(X_1) \sim U[0, 1]$. □