



МЕХАНИКО-
МАТЕМАТИЧЕСКИЙ
ФАКУЛЬТЕТ
МГУ ИМЕНИ
М.В. ЛОМОНОСОВА

teach-in
ЛЕКЦИИ УЧЕНЫХ МГУ

МАТЕМАТИЧЕСКАЯ СТАТИСТИКА

ШАБАНОВ
ДМИТРИЙ АЛЕКСАНДРОВИЧ

МЕХМАТ МГУ

КОНСПЕКТ ПОДГОТОВЛЕН
СТУДЕНТАМИ, НЕ ПРОХОДИЛ
ПРОФ. РЕДАКТУРУ И МОЖЕТ
СОДЕРЖАТЬ ОШИБКИ.
СЛЕДИТЕ ЗА ОБНОВЛЕНИЯМИ
НА [VK.COM/TEACHINMSU](https://vk.com/teachinmsu).

ЕСЛИ ВЫ ОБНАРУЖИЛИ
ОШИБКИ ИЛИ ОПЕЧАТКИ,
ТО СООБЩИТЕ ОБ ЭТОМ,
НАПИСАВ СООБЩЕСТВУ
[VK.COM/TEACHINMSU](https://vk.com/teachinmsu).

Содержание

1.	Лекция 1	7
	Введение	7
	Примеры статистических задач	8
	Вероятностно-статистическая модель	9
	Выборка и эмпирическое распределение	9
	Теорема (Гливленко-Кантелли)	12
	Выборка и эмпирическое распределение	14
2.	Лекция 2	18
	Сходимости случайных векторов	18
	Слабая сходимость вероятностных мер	19
	Предельные теоремы	20
	Теорема о наследовании сходимости	21
	Лемма Слуцкого	22
	Асимптотическая нормальность	26
3.	Лекция 3	27
	Статистики и оценки	27
	Свойства оценок	30
	Наследование свойств	32
	Методы построения оценок	33
4.	Лекция 4	37
	Выборочные квантили	37
	Выборочная медиана	39
	Сравнение оценок	40
5.	Лекция 5	45
	Доминируемые семейства	45
	Условия регулярности	46
	Неравенство Рао-Крамера	47
	Критерий эффективности	48
	О выполнении условий регулярности	50
	Информация Фишера статистик	51
6.	Лекция 6	56
	Многомерный вариант	56
	Матричное неравенство Коши-Буняковского	57
	Условия регулярности	58
	Многомерное неравенство Рао-Крамера	58

	Функция правдоподобия	60
7.	Лекция 7	63
	Экстремальное свойство правдоподобия	63
	Состоятельность решения уравнения правдоподобия	65
	Асимптотическая нормальность решения уравнения правдоподобия	67
	Теорема Бахадура	70
	Эффективность оценки максимального правдоподобия	72
8.	Лекция 8	74
	Теорема Бахадура	74
	Вспомогательная теорема	75
	Доказательство теоремы Бахадура	81
9.	Лекция 9	83
	Условное математическое ожидание	83
	Существование УМО	84
	Дискретные σ -алгебры	86
	Свойства УМО	87
10.	Лекция 10	92
	Условное математическое ожидание	92
	Условное распределение и условная плотность	94
	Теорема о вычислении УМО	94
	Вычисление условной плотности	95
	Схема вычисления УМО	96
	Теорема о наилучшем квадратичном прогнозе	98
11.	Лекция 11	100
	Постановка задачи	100
	Байесовские оценки	101
	Минимаксный подход	104
12.	Лекция 12	109
	Понятие оптимальной оценки	109
	Достаточные статистики	109
	Полные статистики	112
	Критерий факторизации	113
	Примеры	116
13.	Лекция 13	119
	Экспоненциальные семейства	119

Плотность $S(\mathbf{X})$	120
Достаточные статистики и информация Фишера	121
Преобразование Лапласа	122
Доказательство теоремы об экспоненциальном семействе	123
Пример нахождения оптимальной оценки	123
Доверительное оценивание	125
Метод центральной статистики	126
Асимптотические доверительные интервалы	127
14. Лекция 14	130
Постановка задачи линейной регрессии	130
Метод наименьших квадратов	131
Гауссовская линейная модель	134
Доверительные интервалы и области	135
15. Лекция 15	140
Гипотезы	140
Статистические критерии	141
Сравнение критериев	143
Лемма Неймана-Пирсона	144
Монотонное отношение правдоподобия	147
16. Лекция 16	151
Примеры применения теоремы о монотонном отношении правдоподобия	151
Двойственность проверки гипотез и доверительного оценивания	154
Линейные гипотезы	155
Построение критерия	156
Обобщенный метод наименьших квадратов	158
17. Лекция 17	162
Критерии согласия	162
Критерий хи-квадрат	163
Параметрический критерий хи-квадрат	166
Критерий независимости хи-квадрат	168
Критерий однородности хи-квадрат	170
Критерии для непрерывных распределений	171
18. Лекция 18	174
Напоминание из курса теории вероятностей	174
Случайные процессы	175
Измеримые отображения	176

Сходимость случайных процессов	178
Принцип инвариантности	180
19. Лекция 19	182
Теорема Колмогорова	182
20. Лекция 20	192
Мотивировка	192
Последовательные критерии	193
Байесовский подход к проверке гипотез	194
Оптимальность последовательного критерия	199
Приближенные вычисления	200

1. Лекция 1

Введение

Предмет изучения математической статистики

Разберемся с тем, что изучает математическая статистика, какие задачи ставятся и какие цели мы хотим достигать.

Начнем с того, что является предметом изучения **теории вероятностей**. Кратко его можно охарактеризовать как математический анализ случайных явлений. Если говорить неформально, то у нас есть природа некоторого явления (с математической точки зрения это значит, что мы знаем распределение), и мы хотим выяснить как ведут себя характеристики, наблюдаемые в эксперименте.

Математическая статистика - это часть теории вероятностей, в которой изучаются обратные задачи, то есть нам известны экспериментальные данные, и требуется вынести суждение о природе случайного явления.

Пример

Рассмотрим классический пример, который показывает разницу в постановках задач математической статистики и теории вероятностей.

Пусть в городе имеется всего N жителей, среди них есть M заболевших *COVID* – 19.

Задача теории вероятностей

Какова вероятность того, что в случайной выборке из n человек будет ровно m заболевших?

В этой задаче M известно (без него задачу решить не получится). Ответом в этой задаче является гипергеометрическое распределение.

В задаче математической статистики все наоборот. Мы провели наблюдение и выяснили, что среди n жителей ровно m заболевших.

Задача математической статистики

Среди случайно выбранных n жителей оказалось ровно m заболевших. Как можно оценить общее число заболевших в городе?

Здесь M уже выступает в роли неизвестного параметра.

Примеры статистических задач

Рассмотрим модель простой линейной регрессии:

$$y_i = \theta \cdot x_i + \varepsilon_i, \quad i = 1, \dots, n.$$

Здесь y_i, x_i - известные величины, θ - неизвестный параметр, который мы хотели бы оценить, ε_i - случайные неизвестные ошибки измерений. Это простая модель, более общую мы рассмотрим позже.

В качестве **примера** рассмотрим движение объекта по прямой. У нас есть прибор, которым мы измеряем движение объекта, но мы не знаем его скорость. Мы предполагаем, что прибор неточный, поэтому есть некоторые ошибки измерений. В итоге мы получаем $\theta \cdot x_i$ - настоящее положение объекта, с некоторой ошибкой ε_i .

- x_i - моменты измерения положения объекта,
- y_i - измеренное положение объекта в момент времени x_i ,
- ε_i - случайная ошибка измерений,
- θ - неизвестная скорость объекта.

Можно выделить следующие постановки статистических задач:

1) **Точечное оценивание.**

Здесь мы хотим отыскать такую функцию $\hat{\theta} = f(x_1, \dots, x_n, y_1, \dots, y_n)$, что в некотором смысле приближает неизвестную нам θ .

2) **Интервальное оценивание.**

Здесь нам не так важно знать точное значение θ . Мы хотим получить это значение с некоторой точностью, например, указав такое $\varepsilon > 0$, что $|\hat{\theta} - \theta| < \varepsilon$ с большой вероятностью, например, с вероятностью 0,95.

3) **Проверка статистических гипотез**

Здесь требуется проверить (подтвердить или опровергнуть) на основе данных $x_1, \dots, x_n, y_1, \dots, y_n$ некоторое суждение вида $\theta \geq \theta_0$ или $\theta = \theta_0$. Например, нас интересует, произойдет ли повышение уровня воды в реке выше какой-либо критической отметки, причем неважно насколько выше.

4) **Проверка однородности**

Рассмотрим две группы данных X_1, \dots, X_n и Y_1, \dots, Y_m , полученных в результате проведения случайного эксперимента в разных условиях, например, измерения разными приборами.

Необходимо выяснить, верно ли, что условия не влияют на результат? То есть верно ли, что $X_i \stackrel{d}{=} Y_j$ (данные одинаково распределены).

5) Проверка независимости

Пусть результат эксперимента имеет два параметра (фактора) (A, B) , принимающие конечное число значений, $a_1, \dots, a_k, b_1, \dots, b_m$. После проведения серии независимых экспериментов получены количественные данные $n_{ij}, i = 1, \dots, k, j = 1, \dots, m$, по сочетаниям значений факторов.

Верно ли, что факторы A и B независимы?

Математическая статистика - это теория принятия оптимальных статистических решений. Цель - найти оптимальные решения, основанные на статистических данных.

Вероятностно-статистическая модель

Определение 1.1. Пусть имеется наблюдение \mathbf{X} (результат проведенного эксперимента), которое мы считаем случайным элементом (случайным вектором и т.д.). Множество всех возможных значений \mathbf{X} называется **выборочным пространством** и обозначается \mathcal{X} . Мы считаем, что \mathbf{X} получен как результат случайного выбора из \mathcal{X} с неизвестным распределением P .

Замечание

Пусть P и неизвестно, зачастую нам известен класс распределений \mathcal{P} , которому принадлежит P . Например, мы пронаблюдали некоторые данные, и мы предполагаем, что они взяты из нормального распределения, но нам неизвестны параметры этого распределения.

Если у нас есть такой класс распределений, то мы можем образовать тройку.

Определение 1.2. Тройка $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, \mathcal{P})$, где \mathcal{X} - выборочное пространство, $\mathcal{B}_{\mathcal{X}}$ - σ -алгебра на нем, а \mathcal{P} - семейство распределений (вероятностных мер) на $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$, называется **вероятностно-статистической моделью**.

Как правило $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}) = (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$

Выборка и эмпирическое распределение

Определение 1.3. Если наблюдение $\mathbf{X} = (X_1, \dots, X_n)$ есть набор независимых одинаково распределенных случайных величин, то \mathbf{X} называется **выборкой** размера n из некоторого распределения.

Предположим, что у нас есть выборка $\mathbf{X} = (X_1, \dots, X_n)$ размера n из неизвестного распределения P_X на $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Пусть $B \in \mathcal{B}(\mathbb{R})$. Как можно восстановить $P_X(B)$?

Определение 1.4. Для любого $B \in \mathcal{B}(\mathbb{R})$ положим $P_n^*(B) = \frac{\mu(B)}{n}$, где $\mu(B)$ - это число элементов выборки, попавших в множество B , то есть

$$P_n^*(B) = \frac{1}{n} \sum_{i=1}^n I\{X_i \in B\}.$$

Распределением P_n^* , как функция от $B \in \mathcal{B}(\mathbb{R})$, называется **эмпирическим распределением**, построенным при выборке X_1, \dots, X_n . Это дискретное распределение, оно имеет практически равномерное распределение на элементах X_1, \dots, X_n .

Эмпирическая функция распределения

Насколько эмпирическое распределение близко к настоящему?

Утверждение

Пусть $(X_n, n \in \mathbb{N})$ - набор н.о.р.с.в. (независимых одинаково распределенных величин) с распределением P_X . Тогда для любого $B \in \mathcal{B}(\mathbb{R})$ выполнено

$$P_n^*(B) \xrightarrow{\text{п.н.}} P_X(B), \quad n \rightarrow \infty,$$

где п.н. - почти наверное.

Доказательство.

Заметим, что $P_n^*(B) = \frac{1}{n} \sum_{i=1}^n I\{X_i \in B\}$ - сумма н.о.р.с.в. Согласно УЗБЧ (усиленно-му закону больших чисел) имеем

$$P_n^*(B) \xrightarrow{\text{п.н.}} EI\{X_1 \in B\} = P(X_1 \in B) = P_X(B).$$

То есть наша сумма сходится почти наверное к своему математическому ожиданию $EI\{X_1 \in B\}$. Это то же самое, что вероятность того, что X_1 попал в B , это и есть $P_X(B)$. ■

Замечание

К сожалению, данная сходимостъ является точечной, то есть при большом n мы можем хорошо аппроксимировать $P_X(B)$ для конкретного B , но аппроксимация настоящего распределения для всех множеств одновременно невозможна. Однако если рассмотреть "хорошие" множества B , то равномерная аппроксимация получается.

А можно ли аппроксимировать все множества одновременно? Ответ на этот вопрос строго отрицательный, даже для самых простых распределений.

Например, предположим, что у нас есть выборка из равномерного распределения на интервале $(0, 1)$, то есть $X_1, \dots, X_n \sim U(0, 1)$. Мы хотим выяснить чему равно выражение $\sup_{B \in \mathcal{B}(\mathbb{R})} |P_n^*(B) - P_X(B)|$. Это случайная величина, к тому же не совсем понятно, насколько она измерима, так как мы берем супремум по всей борелевской σ -алгебре. На самом деле, это выражение почти всегда равно 1, так как эмпирическое распределение - это дискретное распределение, сосредоточенное в числах X_1, \dots, X_n . Если мы возьмем множество $B = [0, 1] \setminus \{X_1, \dots, X_n\}$, то $P_n^* = 0$, так как там нет элементов выборки, соответственно его вероятность с точки зрения эмпирического распределения нулевая, а настоящее распределение $P_X(B) = 1$, так как мы выкинули из отрезка конечное число точек, следовательно мера B совпадает с мерой отрезка. Получается, что постановка задачи такого вида бессмысленна, то есть мы не можем приближать одновременно все множества.

Чтобы рассмотреть нечто похожее, необходимо взять супремум не по борелевской σ -алгебре, а по множеству из некоторого класса.

Мы помним из теории вероятностей, что само распределение как вероятностная мера однозначно определяется своей функцией распределения. Соответственно, у эмпирического распределения тоже есть функция распределения.

Эмпирическая функция распределения

Определение 1.5. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ есть выборка размера n из неизвестного распределения P_X . Тогда величина

$$F_n^*(x) = P_n^*((-\infty, x]) = \frac{1}{n} \sum_{i=1}^n \mathbf{I}\{X_i \leq x\}$$

называется **эмпирической функцией распределения** выборки X_1, \dots, X_n .

Попробуем нарисовать ее график. Эмпирическое распределение является дискретным, соответственно мы имеем дело с функцией распределения дискретного распределения.

Для удобства предположим, что выборка упорядочена, то есть $X_1 < X_2 < \dots < X_n$. Тогда наша эмпирическая функция распределения выглядит следующим образом:

- До точки X_1 элементов выборки нет, поэтому значение $F_n^*(x)$ равно нулю,
- Когда мы достигли точки X_1 , появляется одна точка, поэтому мы получаем скачок в значение $\frac{1}{n}$ (сумма становится равной 1), и так далее,

- Если элементы выборки совпали, то скачки могут быть на большее значение,
- Синим цветом на графике изобразим настоящую функцию распределения.

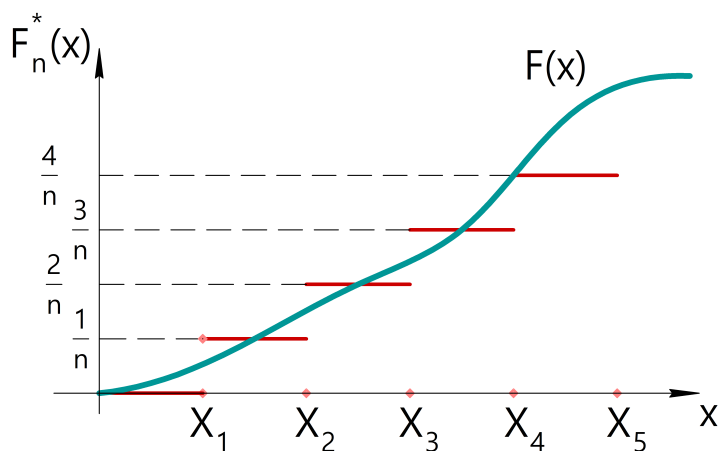


Рис. 1

Возникает вопрос: как хорошо $F_n^*(x)$ приближает настоящую функцию $F(x)$? Мы знаем, что в каждой точке есть сходимости (из утверждения, доказанного выше). Оказывается, что для функции распределения ситуация проще, чем для всего эмпирического распределения, а именно, мы покажем, что $F_n^*(x)$ равномерно приближает $F(x)$ на всей прямой.

Теорема (Гливленко-Кантелли)

Теорема 1.1. (Гливленко-Кантелли)

Пусть $(X_n, n \in \mathbb{N})$ - н.о.р.с.в. с функцией распределения $F(x)$. Тогда

$$\sup_{x \in \mathbb{R}} |F_n^*(x) - F(x)| \xrightarrow{п.н.} 0, \quad b \rightarrow \infty.$$

Доказательство.

Будем считать, что случайные величины $(X_n, n \in \mathbb{N})$ заданы на вероятностном пространстве (Ω, \mathcal{F}, P) . Проверим сначала, что величина

$$D_n(\omega) = \sup_{x \in \mathbb{R}} |F_n^*(x, \omega) - F(x)|$$

является случайной величиной. Заметим, что для каждого ω (элементарный исход) функция $|F_n^*(x, \omega) - F(x)|$ непрерывна справа (это следует из свойств функций распределения) и не может достигать своего максимума на бесконечностях, так как обе функции на $+\infty$ стремятся к 1, а на $-\infty$ к 0. Стало быть,

$$D_n(\omega) = \sup_{x \in \mathbb{Q}} |F_n^*(x, \omega) - F(x)|$$

является супремумом счетного числа ограниченных случайных величин. Значит, D_n - это тоже случайная величина.

Перейдем к доказательству того, что последовательность сходится к нулю п.н. Зафиксируем произвольным образом $N \in \mathbb{N}$ и положим

$$x_{k,N} = \min \left\{ x \in \mathbb{R} : F(x) \geq \frac{k}{N} \right\},$$

$k = 1, \dots, N-1$ и положим $x_{0,N} = -\infty$, $x_{N,N} = +\infty$ (разрежем нашу прямую на отрезки по точкам $x_{k,N}$).

Возьмем произвольный x . Если $x \in [x_{k,N}, x_{k+1,N})$, то попробуем оценить сверху выражение $F_n^*(x) - F(x)$. Воспользуемся свойствами функций распределения, а именно: неубыванием и непрерывностью справа. Тогда:

$$\begin{aligned} F_n^*(x) - F(x) &\leq F_n^*(x_{k+1,N} - 0) - F(x_{k,N}) = \\ &= F_n^*(x_{k+1,N} - 0) - F(x_{k+1,N} - 0) + F(x_{k+1,N} - 0) - F(x_{k,N}) \leq \\ &\leq F_n^*(x_{k+1,N} - 0) - F(x_{k+1,N} - 0) + \frac{1}{N}, \end{aligned}$$

где $F_n^*(x_{k+1,N} - 0)$ - значение предела слева функции F_n^* в точке $x_{k+1,N}$. Значение функции $F(x)$ берется в точке $x_{k,N}$, так как $F(x) \geq F(x_{k,N})$.

$F(x_{k+1,N} - 0) - F(x_{k,N}) \leq \frac{1}{N}$, так как в точке $x_{k,N}$ значение функции хотя-бы $\frac{k}{N}$, а в точке $x_{k+1,N}$ значение хотя-бы $\frac{k+1}{N}$, но это наименьшая точка с таким свойством, поэтому предел слева будет меньше либо равен $\frac{k+1}{N}$.

Аналогично получаем оценку снизу:

$$\begin{aligned} F_n^*(x) - F(x) &\geq F_n^*(x_{k,N}) - F(x_{k+1,N} - 0) = \\ &= F_n^*(x_{k,N}) - F(x_{k,N}) + F(x_{k,N}) - F(x_{k+1,N} - 0) \geq \\ &\geq F_n^*(x_{k,N}) - F(x_{k,N}) - \frac{1}{N}. \end{aligned}$$

Попробуем, с учетом полученного, оценить разность $|F_n^*(x) - F(x)|$ для всех $x \in \mathbb{R}$:

$$\begin{aligned} |F_n^*(x) - F(x)| &\leq \\ &\leq \max_{\substack{1 \leq k \leq N-1 \\ 1 \leq l \leq N-1}} (|F_n^*(x_{k,N}) - F(x_{k,N})|, |F_n^*(x_{l,N} - 0) - F(x_{l,N} - 0)|) + \frac{1}{N}. \end{aligned}$$

Вспомним, что $F_n^*(y) = \frac{1}{n} \sum_{i=1}^n I\{X_i \leq y\}$, а $F_n^*(y-0)$ - предел слева функции распределения, который либо на единицу меньше (оказались в точке X_i), либо равен значению (по непрерывности). Мы посчитаем количество точек, которые строго меньше чем y , то есть $F_n^*(y-0) = \frac{1}{n} \sum_{i=1}^n I\{X_i < y\}$. Согласно УЗБЧ для каждого $y \in \mathbb{R}$ выполнено $F_n^*(y) \xrightarrow{п.н.} F(y)$ и $F_n^*(y-0) \xrightarrow{п.н.} F(y-0)$. Обозначим:

$$\Omega_N = \{\omega : \forall k, l, F_n^*(x_{k,N}, \omega) \longrightarrow F(x_{k,N}), F_n^*(x_{l,N} - 0, \omega) \longrightarrow F(x_{l,N} - 0)\}.$$

Тогда $P(\Omega_N) = 1$ для любого N . В то же время для любого $\omega \in \Omega_N$ выполнено

$$\overline{\lim}_{n \rightarrow \infty} D_n(\omega) \leq \frac{1}{N}.$$

Обозначим $\Omega' = \bigcap_{N=1}^{\infty} \Omega_N$. Тогда $P(\Omega') = 1$ и для любого $\omega \in \Omega'$

$$\exists \overline{\lim}_{n \rightarrow \infty} D_n(\omega) = 0,$$

что и доказывает искомую сходимость $D_n \xrightarrow{п.н.} 0$ при $n \rightarrow +\infty$. ■

Замечание

Теорема Гливенко-Кантелли показывает, что основная задача математической статистики (восстановление неизвестного распределения) вполне обоснована. Например, мы можем равномерно по всем точкам на прямой восстановить функцию распределения.

Глядя на эту теорему возникает вопрос о скорости сходимости D_n к нулю: сколько наблюдений необходимо взять, чтобы можно было сказать, что наша функция распределения эмпирическая? Ответ на этот вопрос очень нетривиален, и получим мы его ближе к концу курса.

Параметрическая модель

Обсудим теперь в каких вероятностных пространствах мы будем работать.

Определение 1.6. Вероятностно-статистическая модель $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, \mathcal{P})$ называется **параметрической**, если семейство \mathcal{P} параметризовано, то есть $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$.

При этом $P_{\theta_1} \neq P_{\theta_2}$, если $\theta_1 \neq \theta_2$.

Примеры

- Биномиальное распределение: $\text{Bin}(1, \theta)$, $\theta \in (0, 1)$.
- Нормальное распределение: $\mathcal{N}(a, \sigma^2)$, $\theta = (a, \sigma^2)$, $a \in \mathbb{R}$, $\sigma > 0$.
- Экспоненциальное распределение: $\text{Exp}(\theta)$, $\theta > 0$.

Преимущество параметрических моделей в том, что если наше распределение параметризовано, то вместо того, чтобы восстанавливать все распределение, мы можем попытаться восстановить конкретное значение параметра, что гораздо проще.

Наблюдение \mathbf{X} мы считаем результатом случайного выбора элемента из \mathcal{X} . Формально, можно понимать \mathbf{X} как тождественное отображение из \mathcal{X} в \mathcal{X} :

$$\mathbf{X}(x) = x \quad \forall x \in \mathcal{X}.$$

Тогда для любого $\theta \in \Theta$ случайный вектор \mathbf{X} определен на вероятностном пространстве $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, P_{\theta})$ и принимает значения в \mathcal{X} .

Замечание

Тогда P_{θ} будет и распределением \mathbf{X} : для любого $B \in \mathcal{B}_{\mathcal{X}}$

$$P_{\theta}(\mathbf{X} \in B) = P_{\theta}(B).$$

Если $\mathbf{X} = (X_1, \dots, X_n)$ - это выборка, то естественно оценивать и параметризовывать распределение P_{θ} только одного элемента выборки. Как тогда построить формальную модель?

Случай конечной выборки

Предположим, что мы построили вероятностно-статистическую модель $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, P)$ для одного элемента выборки. Тогда рассмотрим тройку $(\mathcal{X}^n, \mathcal{B}_{\mathcal{X}}^n, P^n)$, где

- $\mathcal{X}^n = \underbrace{\mathcal{X} \times \dots \times \mathcal{X}}_{n \text{ штук}};$
- $\mathcal{B}_{\mathcal{X}}^n = \underbrace{\mathcal{B}_{\mathcal{X}} \otimes \dots \otimes \mathcal{B}_{\mathcal{X}}}_{n \text{ штук}}$ - σ -алгебра, порожденная прямоугольниками.

Формально

$$\mathcal{B}_{\mathcal{X}}^n = \sigma(B_1 \times \dots \times B_N : B_i \in \mathcal{B}_{\mathcal{X}}).$$

- $P^n = \{P_{\theta}^n, \theta \in \Theta\}$, где $P_{\theta}^n = P_{\theta} \otimes \dots \otimes P_{\theta}$ - вероятностная мера на $(\mathcal{X}^n, \mathcal{B}_{\mathcal{X}}^n)$, которая на прямоугольниках задается по правилу

$$P_{\theta}^n(B_1 \times \dots \times B_n) = \prod_{i=1}^n P_{\theta}(B_i).$$

Лемма

Вероятностная мера P_θ^n существует и единственна с подобным свойством.

Замечание

Говорят, что вероятностное пространство $(\mathcal{X}^n, \mathcal{B}_{\mathcal{X}}^n, P^n)$ является прямым произведением вероятностных мер $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, P)$.

Как моделировать выборку? Мы хотим, чтобы случайные величины были независимы и одинаково распределены, и хотим иметь одно и то же распределение P_θ . Для каждого $i = 1, \dots, n$ определим отображение $X_i : \mathcal{X}^n \rightarrow \mathcal{X}$ по правилу

$$X_i(x_1, \dots, x_n) = x_i.$$

Тогда для любых $B_1, \dots, B_n \in \mathcal{B}_{\mathcal{X}}$ выполнено

$$P_\theta^n(X_1 \in B_1, \dots, X_n \in B_n) = P_\theta^n(B_1 \times \dots \times B_n) = \prod_{i=1}^n P_\theta(B_i) = \prod_{i=1}^n P_\theta^n(X_i \in B_i)$$

Значит, X_1, \dots, X_n - независимы и одинаково распределены с одним и тем же распределением $P_\theta(B_i)$.

Случай бесконечной выборки

В асимптотических вопросах возникает потребность в выборке $\{X_n, n \in \mathbb{N}\}$ неограниченного размера. В этом случае рассматриваем тройку $(\mathcal{X}^\infty, \mathcal{B}_{\mathcal{X}}^\infty, P^\infty)$, где

- $\mathcal{X}^\infty = \{(x_1, x_2, \dots) : x_i \in \mathcal{X}\}$ - пространство последовательностей из \mathcal{X} ;
- $\mathcal{B}_{\mathcal{X}}^\infty$ - цилиндрическая σ -алгебра на \mathcal{X}^∞ . Пусть $B_n \in \mathcal{B}_{\mathcal{X}}^n$, тогда определим цилиндр с основанием B_n ;

$$\mathcal{F}_n(B_n) = \{x \in \mathcal{X}^\infty, x = (x_1, x_2, \dots) : (x_1, \dots, x_n) \in B_n\}$$

Тогда $\mathcal{B}_{\mathcal{X}}^\infty$ - σ -алгебра, порожденная всеми цилиндрами:

$$\mathcal{B}_{\mathcal{X}}^\infty = \sigma(\mathcal{F}_n(B_n) : n \in \mathbb{N}, B_n \in \mathcal{B}_{\mathcal{X}}^n).$$

- $P^\infty = \{P_\theta^\infty, \theta \in \Theta\}$, где P_θ^∞ - вероятностная мера на $(\mathcal{X}^\infty, \mathcal{B}_{\mathcal{X}}^\infty)$, которая на цилиндрах задается по правилу

$$P_\theta^\infty(\mathcal{F}_n(B_n)) = P_\theta^n(B_n).$$

Теорема 1.2. (Колмогоров)

Если $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}) = (\mathbb{R}, \mathcal{B}(\mathbb{R}))$, то вероятностная мера P_θ^∞ существует и единственная с подобным свойством.

Как моделировать выборку? Для каждого $i \in \mathbb{N}$ определим отображение $X_i : \mathcal{X}^\infty \rightarrow \mathcal{X}$ по правилу (такому же, как в n -мерном пространстве)

$$X_i(x_1, x_2, \dots) = x_i$$

Тогда для любых $B_1, \dots, B_n \in \mathcal{B}_{\mathcal{X}}$ выполнено

$$\begin{aligned} P_\theta^\infty(X_1 \in B_1, \dots, X_n \in B_n) &= P_\theta^\infty(\mathcal{F}_n(B_1 \times \dots \times B_n)) = \\ &= P_\theta^n(B_1 \times \dots \times B_n) = \prod_{i=1}^n P_\theta(B_i) = \prod_{i=1}^n P_\theta^\infty(X_i \in B_i). \end{aligned}$$

Значит, $(X_n, n \in \mathbb{N})$ - независимы и одинаково распределены с распределением P_θ .

Замечание

В дальнейшем мы будем опускать индексы n и ∞ у мер P_θ^n и P_θ^∞ , обозначая их также P_θ , как и распределение одного наблюдения X_i .

В записях вида

$$P_\theta(\mathbf{X} \in B), \quad E_\theta f(\mathbf{X})$$

предполагается, что вероятности и математические ожидания берутся в пространстве $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, P_\theta)$, то есть при условии, что θ - истинное значение параметра. Отметим, что в построенной модели смена параметра не требует смены отображения \mathbf{X} .

В параметрической модели вопрос о нахождении истинного распределения сводится к вопросу о нахождении истинного значения параметра (что, конечно, выглядит куда проще).

2. Лекция 2

Формально, материал этой лекции относится к теории вероятностей, но в силу того, что зачастую в базовом курсе он отсутствует, мы пройдем его в нашем курсе и изучим те основные инструменты работы со сходимости, которые нам будут крайне нужны в рамках математической статистики.

Сходимости случайных векторов

Определение 2.1. Пусть $\{\xi_n, n \in \mathbb{N}\}$, ξ - случайные векторы размерности m . Тогда последовательность ξ_n сходится к ξ

- с вероятностью 1 (почти наверное, $\xi_n \xrightarrow{\text{п.н.}} \xi$), если

$$P\left(\lim_{n \rightarrow +\infty} \xi_n = \xi\right) = 1.$$

- по вероятности ($\xi_n \xrightarrow{P} \xi$), если $\forall \varepsilon > 0$ выполнено

$$P(\|\xi_n - \xi\|_2 \geq \varepsilon) \rightarrow 0 \quad \text{при } n \rightarrow \infty,$$

где $\|x\|_2 = \sqrt{x_1^2 + \dots + x_m^2}$ для $x \in \mathbb{R}^m$.

- по распределению ($\xi_n \xrightarrow{d} \xi$), если для любой ограниченной непрерывной функции $f: \mathbb{R}^m \rightarrow \mathbb{R}$ выполнено

$$Ef(\xi_n) \rightarrow Ef(\xi) \quad \text{при } n \rightarrow \infty.$$

Естественно задать вопрос: чем отличаются данные сходимости от одномерных?

Упражнение (1)

Для сходимости почти наверное и по вероятности векторная сходимость эквивалентна соответствующим сходимостям компонент: если $\xi = (\xi^{(1)}, \dots, \xi^{(m)})$, $\xi_n = (\xi_n^{(1)}, \dots, \xi_n^{(m)})$, то

$$\xi_n \xrightarrow{\text{п.н.}} \xi \Leftrightarrow \forall i = 1, \dots, m \quad \xi_n^{(i)} \xrightarrow{\text{п.н.}} \xi^{(i)},$$

$$\xi_n \xrightarrow{P} \xi \Leftrightarrow \forall i = 1, \dots, m \quad \xi_n^{(i)} \xrightarrow{P} \xi^{(i)}.$$

Упражнение (2)

Для сходимости по распределению векторная сходимость влечет сходимость компонент: если $\xi = (\xi^{(1)}, \dots, \xi^{(m)})$, $\xi_n = (\xi_n^{(1)}, \dots, \xi_n^{(m)})$, то

$$\xi_n \xrightarrow{d} \xi \Leftrightarrow \forall i = 1, \dots, m \quad \xi_n^{(i)} \xrightarrow{d} \xi^{(i)}.$$

Обратное утверждение, вообще говоря, неверно.

Упражнение (3, взаимоотношение видов сходимостей)

$$\xi_n \xrightarrow{\text{п.н.}} \xi \Rightarrow \xi_n \xrightarrow{P} \xi \Rightarrow \xi_n \xrightarrow{d} \xi.$$

Слабая сходимость вероятностных мер

Слабая сходимость вероятностных мер пригодится нам в дальнейшем, когда мы будем говорить о доказательстве критерия Колмогорова. Пока напомним о том, что это такое и поговорим о том, как интерпретируется сходимость по распределению в терминах сходимости самих распределений.

Пусть (S, ρ) - метрическое пространство.

Определение 2.2. Борелевской σ -алгеброй, $\mathcal{B}(S)$, на (S, ρ) называется минимальная σ -алгебра, содержащая все открытые множества в S .

Определение 2.3. Пусть задано метрическое пространство (S, ρ) и последовательность $\{Q_n, n \in \mathbb{N}\}$ вероятностных мер на $(S, \mathcal{B}(S))$. Будем говорить, что Q_n слабо сходится к вероятностной мере Q на $(S, \mathcal{B}(S))$, если для любой ограниченной непрерывной функции $f: S \rightarrow \mathbb{R}$ выполнено

$$\lim_{n \rightarrow \infty} \int_S f(x) Q_n(dx) = \int_S f(x) Q(dx).$$

Обозначение: $Q_n \xrightarrow{w} Q$.

Теорема Александрова

Утверждение Если пространство (S, ρ) сепарабельно, то $\mathcal{B}(S)$ является минимальной σ -алгеброй, содержащей все открытые шары.

Одним из главных инструментов работы со слабой сходимостью является следующая теорема Александрова, которая дает серию эквивалентных утверждений, каждое из которых эквивалентно слабой сходимости. Соответственно, чтобы доказать слабую сходимость, будет достаточно проверить одно из свойств, представленных в теореме.

Теорема 2.1. (А.Д. Александров)

Пусть $\{Q_n, n \in \mathbb{N}\}$ и Q - вероятностные меры на метрическом пространстве (S, ρ) . Тогда следующие утверждения эквивалентны:

1) $Q_n \xrightarrow{w} Q$,

- 2) $\overline{\lim}_{n \rightarrow \infty} Q_n(F) \leq Q(F)$ для любого замкнутого множества $F \subset S$,
- 3) $\underline{\lim}_{n \rightarrow \infty} Q_n(G) \geq Q(G)$ для любого открытого множества $G \subset S$,
- 4) Для любого борелевского множества $B \in \mathcal{B}(S)$ такого, что $Q(\partial B) = 0$, выполнено $Q_n(B) \rightarrow Q(B)$ при $n \rightarrow \infty$.

Пояснение

Пусть есть некоторое множество $B \subset S$. Рассмотрим его замыкание $[B]$ в S (для этого необходимо дополнить множество B всеми предельными точками). Тогда граница $\partial B = [B] \cap [\overline{B}]$, где \overline{B} - дополнение к множеству B .

Обозначение для свойства 4: $Q_n \Rightarrow Q$ (сходимость в основном).

Доказательство теоремы Александрова можно найти в курсе теории вероятностей.

Взаимоотношение видов сходимостей

Замечание

Пусть P_η обозначает распределение случайного вектора η . Тогда

$$\xi_n \xrightarrow{d} \xi \Leftrightarrow P_{\xi_n} \xrightarrow{w} P_\xi \Leftrightarrow (\text{т. Александрова}) P_{\xi_n} \Rightarrow P_\xi.$$

Отсюда можно вывести факт, который объясняет, зачем вообще нужна сходимость по распределению. Рассмотрим теорему, которая говорит о том, чему эквивалентна сходимость по распределению для случайных величин.

Теорема 2.2. Если $\{\xi_n, n \in \mathbb{N}\}$, ξ - случайные величины, то $\xi_n \xrightarrow{d} \xi$ тогда и только тогда, когда $F_{\xi_n}(x) \rightarrow F_\xi(x)$ для всех $x \in \mathbb{C}(F)$, где $\mathbb{C}(F)$ - множество точек непрерывности функции распределения $F_\xi(x)$.

Для случайных векторов мы можем тоже предложить подобное утверждение.

Предельные теоремы

Вспомним предельные теоремы, и где вообще могут возникать сходимости.

Утверждение (УЗБЧ для случайных векторов)

Пусть $\{X_n, n \in \mathbb{N}\}$ - последовательность независимых одинаково распределенных случайных векторов из \mathbb{R}^m , EX_1 конечно. Тогда

$$\frac{X_1 + \dots + X_n}{n} \xrightarrow{\text{п.н.}} EX_1.$$

Доказательство.

Следует из одномерного УЗБЧ и эквивалентности векторной сходимости п.н. сходимости п.н. всех компонент. ■

Теорема 2.3. (многомерная центральная предельная теорема)

Пусть $\{X_n, n \in \mathbb{N}\}$ - независимые одинаково распределенные случайные векторы из \mathbb{R}^m , $EX_n = a$, $DX_n = \Sigma$ (матрица ковариации). Обозначим $S_n = X_1 + \dots + X_n$. Тогда

$$\sqrt{n} \left(\frac{S_n}{n} - a \right) \xrightarrow{d} \mathcal{N}(0, \Sigma).$$

Теорема о наследовании сходимости

Ключевым утверждением, позволяющим работать со сходимостями, является теорема о наследовании сходимости.

Теорема 2.4. (о наследовании сходимости)

Пусть $\{\xi_n, n \in \mathbb{N}\}$, ξ - случайные векторы размерности m . Пусть $h(x) : \mathbb{R}^m \rightarrow \mathbb{R}^k$ - функция, непрерывная почти всюду относительно распределения ξ (т.е. $\exists B \in \mathcal{B}(\mathbb{R}^m)$ такое, что h непрерывна на B и $P(\xi \in B) = 1$). Тогда

- 1) $\xi_n \xrightarrow{п.н.} \xi \Rightarrow h(\xi_n) \xrightarrow{п.н.} h(\xi),$
- 2) $\xi_n \xrightarrow{P} \xi \Rightarrow h(\xi_n) \xrightarrow{P} h(\xi),$
- 3) $\xi_n \xrightarrow{d} \xi \Rightarrow h(\xi_n) \xrightarrow{d} h(\xi).$

Доказательство первого и второго пункта довольно простые. В третьем пункте придется воспользоваться теоремой Александрова.

Если бы h была просто непрерывной функцией, то доказательство стало бы намного проще.

Покажем это для третьего пункта: необходимо показать, что $h(\xi_n) \xrightarrow{d} h(\xi)$, где $h(\xi_n) \in \mathbb{R}^k$. Возьмем ограниченную непрерывную функцию $f : \mathbb{R}^k \rightarrow \mathbb{R}$. Необходимо удостовериться в том, что $Ef(h(\xi_n)) \rightarrow Ef(h(\xi))$. Так как функция $f \circ h : \mathbb{R}^m \rightarrow \mathbb{R}$ тоже ограничена и непрерывна, так как композиция непрерывных функций - непрерывная функция, а внешняя функция ограничена, то $Ef(h(\xi_n)) \rightarrow Ef(h(\xi))$ верно в силу того, что $\xi_n \xrightarrow{d} \xi$. Перейдем к доказательству теоремы.

Доказательство.

1.

$$P(\lim_{n \rightarrow \infty} h(\xi_n) = h(\xi)) \geq P(\lim_{n \rightarrow \infty} h(\xi_n) = h(\xi), \xi \in B) \stackrel{\substack{\text{т.к. } h \text{ непр.} \\ \text{на } B}}{\geq} P(\lim_{n \rightarrow \infty} \xi_n = \xi, \xi \in B) = 1$$

так как оба события имеют полную вероятность.

2. Пусть $h(\xi_n) \not\xrightarrow{P} h(\xi)$. Тогда $\exists \varepsilon_0, \delta_0 > 0$ и подпоследовательность ξ_{n_k} такая, что

$$P(\|h(\xi_{n_k}) - h(\xi)\|_2 \geq \varepsilon_0) \geq \delta_0 \quad \forall k.$$

Но $\xi_{n_k} \xrightarrow{P} \xi$, значит, существует подпоследовательность, сходящаяся почти наверное: $\xi_{n_{k_s}} \xrightarrow{п.н.} \xi$ при $s \rightarrow \infty$. Согласно 1. получаем, что $h(\xi_{n_{k_s}}) \xrightarrow{п.н.} h(\xi)$.

Значит, $h(\xi_{n_{k_s}}) \xrightarrow{P} h(\xi)$ при $s \rightarrow \infty$. Противоречие с выбором подпоследовательности ξ_{n_k} .

3. Обозначим Q_n - распределение $h(\xi_n)$, Q - распределение $h(\xi)$. Хотим показать, что $Q_n \xrightarrow{w} Q$. По теореме Александрова достаточно показать, что для любого замкнутого $F \subset \mathbb{R}^k$ выполнено

$$\lim_n Q_n(F) \leq Q(F).$$

Имеем,

$$\lim_n Q_n(F) = \lim_n P_{\xi_n}(h^{-1}(F)) \leq \lim_n P_{\xi_n}([h^{-1}(F)]) \stackrel{\text{т.к. } P_{\xi_n} \xrightarrow{w} P_{\xi}}{\leq} P_{\xi}([h^{-1}(F)]).$$

Но в силу замкнутости F

$$[h^{-1}(F)] \subset \overline{B} \cup h^{-1}(F),$$

ведь если $x_n \rightarrow x$, $x_n \in h^{-1}(F)$ и $x \in B$, то $h(x) \in F$. Учитывая, что $P_{\xi}(\overline{B}) = 0$, получаем, что $P_{\xi}([h^{-1}(F)]) = P_{\xi}(h^{-1}(F)) = Q(F)$. ■

Лемма Слуцкого

Второй важнейший инструмент работы со сходимостями - это лемма Слуцкого.

Теорема 2.5. (лемма Слуцкого)

Пусть $\{\xi, n \in \mathbb{N}\}$ и $\{\eta, n \in \mathbb{N}\}$ - две последовательности случайных векторов (вообще говоря, разной размерности). Если $\xi_n \xrightarrow{d} \xi$ и $\eta_n \xrightarrow{d} C$, где C - константа, тогда

$$(\xi_n, \eta_n) \xrightarrow{d} (\xi, C).$$

Мы упоминали в начале лекции, что сходимость компонент по распределению не означает сходимости векторной по распределению. Это крайне неудобно, так как

теорема о наследовании сходимости требует векторной сходимости. Утверждение леммы Слуцкого состоит в том, что если одна из координат сходится по распределению к константе, то выполнена и двумерная сходимость по распределению. В то же время, теорема о наследовании сходимости позволяет брать функции от многомерных сходимостей. Тем самым, можно удачно сочетать данные два инструмента.

Доказательство.

Пусть $\xi \in \mathbb{R}^m$, $C \in \mathbb{R}^k$. Необходимо показать, что для любой ограниченной непрерывной функции $f : \mathbb{R}^{m+k} \rightarrow \mathbb{R}$ выполнено

$$Ef(\xi_n, \eta_n) \rightarrow Ef(\xi, C) \quad \text{при } n \rightarrow \infty$$

Заметим, что достаточно проверить указанную сходимость только для равномерно непрерывных функций f . Например, для сходимости характеристических функций нам даже достаточно полагать, что $f(x) = \cos\langle x, t \rangle$ или $f(x) = \sin\langle x, t \rangle$.

Зафиксируем $\varepsilon > 0$. Положим $M = \max_{x \in \mathbb{R}^{m+k}} |f(x)|$. Выберем $\delta > 0$ так, чтобы для любого $n \geq n_0(\varepsilon)$ были выполнены условия:

•

$$2M \cdot P(\|\eta_n - C\| \geq \delta) \leq \frac{\varepsilon}{2};$$

• для любых $x \in \mathbb{R}^m$ и $\|y - y'\| \leq 2\delta$ верно, что

$$|f(x, y) - f(x, y')| \leq \frac{\varepsilon}{2}.$$

Далее рассмотрим функции

$$g_1(x) = \max_{y: \|y-C\| \leq \delta} f(x, y), \quad g_2(x) = \min_{y: \|y-C\| \leq \delta} f(x, y).$$

Это непрерывные ограниченные функции на \mathbb{R}^m , причем

$$|g_i(x) - f(x, C)| \leq \frac{\varepsilon}{2}.$$

Теперь разобьем искомую величину на сумму

$$Ef(\xi_n, \eta_n) = E(f(\xi_n, \eta_n)I\{\|\eta_n - C\| \leq \delta\}) + E(f(\xi_n, \eta_n)I\{\|\eta_n - C\| > \delta\}).$$

При $n \geq n_0(\varepsilon)$ второе слагаемое не превосходит $M \cdot P(\|\eta_n - C\| \geq \delta) \leq \frac{\varepsilon}{4}$.

Первое же слагаемое можно оценить следующим образом:

$$E(g_2(\xi_n)I\{\|\eta_n - C\| \leq \delta\}) \leq E(f(\xi_n, \eta_n)I\{\|\eta_n - C\| \leq \delta\}) \leq E(g_1(\xi_n)I\{\|\eta_n - C\| \leq \delta\}).$$

Далее, заметим, что

$$\begin{aligned} |Eg_i(\xi_n) - E(g_i(\xi_n)I\{\|\eta_n - C\| \leq \delta\})| &\leq E(|g_i(\xi_n)|I\{\|\eta_n - C\| > \delta\}) \leq \\ &\leq M \cdot P(\|\eta_n - C\| \geq \delta) \leq \frac{\varepsilon}{4}. \end{aligned}$$

Из приведенных неравенств получаем, что при $n \geq n_0(\varepsilon)$

$$Eg_2(\xi_n) - \frac{\varepsilon}{2} \leq Ef(\xi_n, \eta_n) \leq Eg_1(\xi_n) + \frac{\varepsilon}{2}.$$

Перейдем к пределу при $n \rightarrow +\infty$ и используем тот факт, что $\xi_n \xrightarrow{d} \xi$:

$$Eg_2(\xi) - \frac{\varepsilon}{2} \leq \liminf_{n \rightarrow \infty} Ef(\xi_n, \eta_n) \leq \overline{\lim}_{n \rightarrow \infty} Ef(\xi_n, \eta_n) \leq Eg_1(\xi) + \frac{\varepsilon}{2}.$$

Остается вспомнить, что для любого x выполнено $|g_i(x) - f(x, C)| \leq \frac{\varepsilon}{2}$. В итоге:

$$Ef(\xi, C) - \varepsilon \leq \liminf_{n \rightarrow \infty} Ef(\xi_n, \eta_n) \leq \overline{\lim}_{n \rightarrow \infty} Ef(\xi_n, \eta_n) \leq Ef(\xi, C) + \varepsilon.$$

В силу произвольности ε получаем, что существует

$$\lim_{n \rightarrow \infty} Ef(\xi_n, \eta_n) = Ef(\xi, C).$$

■

Полезно следующее следствие из леммы Slutsky.

Следствие

Пусть $\{\xi, n \in \mathbb{N}\}$ и $\{\eta, n \in \mathbb{N}\}$ - две последовательности случайных величин. Если $\xi_n \xrightarrow{d} \xi$ и $\eta_n \xrightarrow{d} C$, где C - константа. Тогда

$$\xi_n + \eta_n \xrightarrow{d} \xi + C, \quad \xi_n \cdot \eta_n \xrightarrow{d} \xi \cdot C.$$

Доказательство.

Согласно лемме Slutsky выполнено $(\xi_n, \eta_n) \xrightarrow{d} (\xi, C)$. Тогда применяя теорему о наследовании сходимости для непрерывных функций суммы и произведения в \mathbb{R}^2 , получаем, что

$$\xi_n + \eta_n \xrightarrow{d} \xi + C, \quad \xi_n \cdot \eta_n \xrightarrow{d} \xi \cdot C.$$

■

Пример применения

Рассмотрим следующее применение теоремы о наследовании сходимости и леммы Слуцкого.

Пример

Пусть $\xi_n \xrightarrow{d} \xi$ - случайные величины, а $h(x)$ - функция, дифференцируемая в точке $a \in \mathbb{R}$. Найдите предел сходимости по распределению у последовательности

$$\frac{h(a + b_n \xi_n) - h(a)}{b_n},$$

где $b_n \rightarrow 0$ - произвольная последовательность положительных чисел.

Решение

Рассмотрим функцию

$$H(x) = \begin{cases} \frac{h(x+a) - h(a)}{x}, & x \neq 0; \\ h'(a), & x = 0. \end{cases}$$

Ясно, что H непрерывна в нуле.

По лемме Слуцкого $b_n \xi_n \xrightarrow{d} 0 \cdot \xi = 0$. Следовательно, по теореме о наследовании сходимости

$$H(b_n \xi_n) \xrightarrow{d} H(0) = h'(a).$$

Снова применяем лемму Слуцкого:

$$\xi_n \cdot H(b_n \xi_n) \xrightarrow{d} \xi \cdot h'(a).$$

Остается заметить, что

$$\xi_n \cdot H(b_n \xi_n) = \frac{h(a + b_n \xi_n) - h(a)}{b_n}.$$

Общий случай

Полезно и следующее обобщение разобранного примера на многомерный случай.

Утверждение

Пусть $\xi_n \xrightarrow{d} \xi$ - случайные векторы из \mathbb{R}^m . Пусть $H(x) : \mathbb{R}^m \rightarrow \mathbb{R}^k$ - вектор-функция, дифференцируемая в точке $a \in \mathbb{R}^m$. Обозначим через

$$H'(a) = \left(\frac{\partial H_i(x)}{\partial x_j}, \quad i = 1, \dots, k, \quad j = 1, \dots, m \right) \Big|_{x=a}$$

матрицу частных производных. Тогда для любой последовательности положительных чисел $b_n \rightarrow 0$ выполнено

$$\frac{H(a + b_n \xi_n) - H(a)}{b_n} \xrightarrow{d} H'(a) \xi.$$

Доказательство оставляется в качестве упражнения.

Асимптотическая нормальность

В математической статистике важную роль играют асимптотически нормальные оценки. Один из основных способов их получения состоит в применении данных утверждений к (многомерной) центральной предельной теореме.

Следствие (дельта-метод)

Пусть $\{X, n \in \mathbb{N}\}$ - независимые одинаково распределенные случайные векторы из \mathbb{R}^m , $EX_n = a$, $DX_n = \Sigma$. Обозначим $S_n = X_1 + \dots + X_n$. Пусть $h : \mathbb{R}^m \rightarrow \mathbb{R}$ - дифференцируемая в точке $a \in \mathbb{R}^m$ функция. Тогда

$$\sqrt{n} \left(h \left(\frac{S_n}{n} \right) - h(a) \right) \xrightarrow{d} \langle h'(a), \mathcal{N}(0, \Sigma) \rangle.$$

Решение

Подберем параметры для предыдущего утверждения. Согласно многомерной ЦПТ

$$\sqrt{n} \left(\frac{S_n}{n} - a \right) \xrightarrow{d} \mathcal{N}(0, \Sigma).$$

Тогда положим:

- $\xi_n = \sqrt{n} \left(\frac{S_n}{n} - a \right),$
- $\xi \sim \mathcal{N}(0, \Sigma),$
- $H(x) = h(x),$
- $a = a,$
- $b_n = \frac{1}{\sqrt{n}}$ (положительная и стремится к 0).

Тогда

$$\frac{H(a + b_n \xi_n) - H(a)}{b_n} = \sqrt{n} \left(h \left(\frac{S_n}{n} \right) - h(a) \right),$$
$$H'(a)\xi = h'(a)^T \xi \sim \langle h'(a), \mathcal{N}(0, \Sigma) \rangle.$$

3. Лекция 3

Статистики и оценки

Пусть $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, \mathcal{P})$ - вероятностно-статистическая модель. Пусть \mathbf{X} - наблюдение в этой модели. Пусть (E, \mathcal{E}) - измеримое пространство.

Определение 3.1. Статистикой $S(\mathbf{X})$ называется измеримая функция от наблюдения \mathbf{X} . Напомним, что отображение $S : \mathcal{X} \rightarrow E$ является измеримым, если для любого множества $B \in \mathcal{E}$ выполнено

$$S^{-1}(B) = \{x : S(x) \in B\} \in \mathcal{B}_{\mathcal{X}}.$$

Важный момент: если модель является параметрической, то отображение S не должно зависеть от параметра!

Пусть $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, \mathcal{P})$, $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$ - параметрическая модель. Пусть \mathbf{X} - наблюдение в этой модели.

Определение 3.2. Если статистика $S(\mathbf{X})$ принимает значения в Θ , то ее можно назвать **оценкой** неизвестного параметра $\theta \in \Theta$. Можно также оценивать функции от параметра, $\tau(\theta)$. В этом случае $S(\mathbf{X})$ должна принимать значения в $\tau(\Theta)$.

Замечание

Мы будем использовать данные два термина почти как эквивалентные.

Примеры статистик и оценок

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка, $X_i \in \mathbb{R}$.

Общая идея построения хороших статистик.

Если Q - распределение вероятностей на прямой, а $G(Q)$ - некоторый функционал, то $G(P_n^*)$ - потенциально хорошая статистика. Здесь P_n^* - эмпирическое распределение, построенное по выборке $\mathbf{X} = (X_1, \dots, X_n)$.

1) Выборочные усреднения.

Пусть $g(x)$ - некоторая борелевская функция \mathbb{R} , а

$$G(Q) = \int_{\mathbb{R}} g(x) Q(dx).$$

Тогда, если подставить в качестве аргумента эмпирическое распределение (равномерное распределение на элементах выборки), то мы по сути получим

математическое ожидание функции $g(x)$ по этому распределению, то есть

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

- среднее значение функции $g(x)$ по выборке. Хорошие примеры:

- $g(x) = x$, тогда $\overline{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n X_i$ - выборочное среднее;
- $g(x) = x^k$, тогда $\overline{\mathbf{X}^k} = \frac{1}{n} \sum_{i=1}^n X_i^k$ - выборочный момент порядка k .

2) Функции от выборочных усреднений.

Пусть $g_1(x), \dots, g_k(x)$ - борелевские функции на \mathbb{R} , а $h : \mathbb{R}^k \rightarrow \mathbb{R}$ - борелевская функция нескольких переменных, тогда можно рассмотреть

$$G(P_n^*) = h\left(\overline{g_1(\mathbf{X})}, \dots, \overline{g_k(\mathbf{X})}\right).$$

Хорошие примеры:

- $h(x, y) = x - y^2$, $g_1(x) = x^2$, $g_2(x) = x$, тогда получаем статистику $S^2 = \overline{\mathbf{X}^2} - (\overline{\mathbf{X}})^2$ - выборочную дисперсию;

Упражнение

Проверить, что

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{\mathbf{X}})^2.$$

- $M_k = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{\mathbf{X}})^k$ - выборочный центральный момент порядка k .

3) Порядковые статистики.

Введем обозначения (упорядочим выборку):

$$\begin{aligned} X_{(1)} &= \min\{X_1, \dots, X_n\}, \\ X_{(2)} &= \min\{(X_1, \dots, X_n) \setminus X_{(1)}\}, \\ &\vdots \\ X_{(n)} &= \max\{X_1, \dots, X_n\}. \end{aligned}$$

Вектор $(X_{(1)}, \dots, X_{(n)})$ называется вариационным рядом.

4) Выборочные квантили.

Определение 3.3. Пусть распределение Q имеет функцию распределения $F(x)$, а $p \in (0, 1)$. Тогда p -квантилью распределения Q называется

$$\zeta_p = \min\{x : F(x) \geq p\}.$$

Если $F(x)$ непрерывна, то для любого $p \in (0, 1)$ будет выполнено $F(\zeta_p) = p$.

Изобразим, как может выглядеть эта величина. Возьмем произвольную функцию распределения и рассмотрим три принципиальных случая.

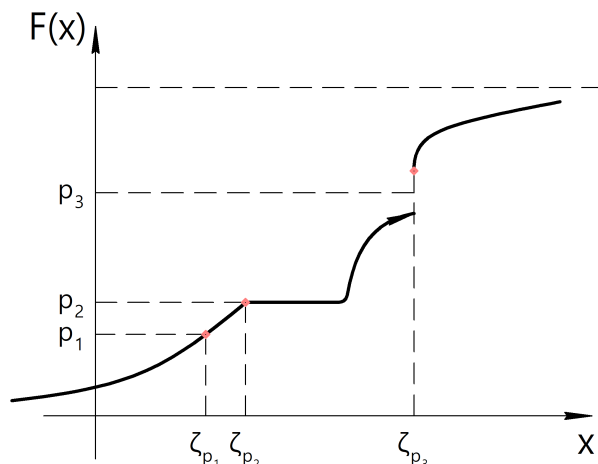


Рис. 2

Первый случай: ζ_{p_1} - квантиль честная, то есть у нас ровно одна точка, в которой достигается значение p_1 .

Второй случай: если у функции распределения есть отрезок постоянства (множество точек, в которых значение равно p_2). Тогда среди этих точек мы берем наименьшую - ζ_{p_2} .

Третий случай: в точке p_3 значение не достигается. Тогда в качестве квантиля мы берем наименьшую точку, в которой значение превышает p_3 - ζ_{p_3} .

Определение 3.4. Выборочной p -квантилью называется p -квантиль эмпирического распределения P_n^* , то есть

$$Z_{n,p} = \begin{cases} X_{([np]+1)}, & \text{если } np \notin \mathbb{Z}; \\ X_{(np)}, & \text{если } np \in \mathbb{Z}. \end{cases}$$

Напомним, что эмпирическое распределение - дискретное, сосредоточенное в точках X_1, \dots, X_n и принимающее значения равновероятно, то есть эмпирическая функция распределения идет "лесенкой" и принимает значения вида $\frac{k}{n}$. Соответственно, почти всегда у нас будет третий случай, когда нет точки, в которой значение нашей функции распределения в точности равно p .

Замечание: в качестве функционала выступает квантиль распределения.

5) М-оценки.

Пусть $\psi(x, \theta)$, $\theta \in \Theta$ - некоторый набор функций. Тогда M -оценкой называется

статистика вида

$$S(X) = \arg \max_{\theta \in \Theta} \left(\frac{1}{n} \sum_{i=1}^n \psi(X_i, \theta) \right).$$

Примером такой оценки будет оценка максимального правдоподобия.

Свойства оценок

Теперь поймем, какие мы хотели бы видеть свойства у оценок, чтобы можно было рассматривать их в качестве хороших.

Несмещённость

Пусть \mathbf{X} - наблюдение с неизвестным распределением $P \in \{P_\theta : \theta \in \Theta\}$, где $\theta \in \mathbb{R}^k$.

Определение 3.5. Оценка $\hat{\theta}(\mathbf{X})$ называется несмещённой оценкой параметра θ , если для любого $\theta \in \Theta$ выполняется равенство

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

Вспомним, что запись E_θ означает, что при взятии математического ожидания мы предполагаем, что выборка \mathbf{X} взята из распределения P_θ :

$$E_\theta \hat{\theta}(\mathbf{X}) = \int_{\mathcal{X}} \hat{\theta}(x) P_\theta(dx).$$

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из $\mathcal{N}(\theta, 1)$. Тогда \bar{X} (выборочное среднее) и X_1 - несмещённые оценки θ .

Состоятельность

В асимптотических вопросах мы предполагаем, что $\mathbf{X} = (X_1, \dots, X_n)$ - это выборка неограниченного размера из неизвестного распределения $P \in \{P_\theta : \theta \in \Theta\}$, $\Theta \in \mathbb{R}^k$.

Определение 3.6. Оценка $\hat{\theta}_n(X_1, \dots, X_n)$ (точнее - последовательность оценок) называется состоятельной оценкой параметра θ , если для любого $\theta \in \Theta$

$$\hat{\theta}(X_1, \dots, X_n) \xrightarrow{P_\theta} \theta, \quad n \rightarrow \infty.$$

Запись означает, что для любого $\theta \in \Theta$ и для любого $\varepsilon > 0$

$$P_\theta \left(\|\hat{\theta}(X_1, \dots, X_n) - \theta\| \geq \varepsilon \right) \rightarrow 0, \quad n \rightarrow \infty.$$

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из $\mathcal{N}(\theta, 1)$. Тогда по ЗБЧ оценка \bar{X} является состоятельной.

Сильная состоятельность

Определение 3.7. Оценка $\hat{\theta}_n(X_1, \dots, X_n)$ (точнее - последовательность оценок) называется сильно состоятельной оценкой параметра θ , если для любого $\theta \in \Theta$

$$\hat{\theta}(X_1, \dots, X_n) \xrightarrow{P_{\theta-\text{п.н.}}} \theta, \quad n \rightarrow \infty.$$

Запись означает, что для любого $\theta \in \Theta$

$$P_{\theta} \left(\lim_{n \rightarrow \infty} \hat{\theta}_n(X_1, \dots, X_n) = \theta \right) = 1.$$

Пример

Пусть $\mathbf{X}(X_1, \dots, X_n)$ - выборка из $\mathcal{N}(\theta, 1)$. Тогда по УЗБЧ оценка \bar{X} является сильно состоятельной.

Асимптотическая нормальность

Ограничимся одномерным вариантом. Пусть $\Theta \subset \mathbb{R}$.

Определение 3.8. Оценка $\hat{\theta}_n(X_1, \dots, X_n)$ (точнее - последовательность оценок) называется асимптотически нормальной оценкой параметра θ , если для любого $\theta \in \Theta$

$$\sqrt{n} \left(\hat{\theta}_n(X_1, \dots, X_n) - \theta \right) \xrightarrow{d_{\theta}} \mathcal{N}(0, \sigma^2(\theta)), \quad n \rightarrow \infty.$$

Функция $\sigma^2(\theta)$ называется асимптотической дисперсией оценки $\hat{\theta}_n$.

Запись означает, что для любого $\theta \in \Theta$ и для любого $x \in \mathbb{R}$ выполнено

$$\lim_{n \rightarrow \infty} P_{\theta} \left(\sqrt{n} \left(\hat{\theta}_n(X_1, \dots, X_n) - \theta \right) \leq x \right) = \Phi \left(\frac{x}{\sigma(\theta)} \right),$$

где Φ - функция Лапласа (функция распределения стандартного нормального закона).

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из $\text{Bin}(1, \theta)$. Тогда по ЦПТ оценка \bar{X} является асимптотически нормальной.

Замечания

- 1) Асимптотические свойства (состоятельность, сильная состоятельность и асимптотическая нормальность) имеют смысл только в случае выборки большого размера. Несмещённость - это "конечное" свойство.
- 2) Оценивать можно не только сам параметр θ , но и функции от него, $\tau(\theta)$. Определения свойств при этом сохраняются.

Наследование свойств

Вопрос: пусть $\hat{\theta}_n(X_1, \dots, X_n)$ хорошо приближает θ , как тогда хорошо приблизить $\tau(\theta)$?

Допустим, что $\hat{\theta}_n$ является состоятельной оценкой. В качестве хорошей оценки для $\tau(\theta)$ можно взять $\tau(\hat{\theta}_n)$:

Утверждение (1)

Пусть $\hat{\theta}_n(X_1, \dots, X_n)$ - сильно состоятельная (состоятельная) оценка параметра θ . Если $\tau(\theta)$ непрерывна на $\Theta \subset \mathbb{R}^k$, то $\tau(\hat{\theta}_n(X_1, \dots, X_n))$ будет сильно состоятельной (состоятельной) оценкой $\tau(\theta)$.

Доказательство. Сразу следует из теоремы о наследовании сходимости. ■

Утверждение (2)

Пусть $\hat{\theta}_n(x_1, \dots, x_n)$ - асимптотически нормальная оценка параметра $\theta \in \mathbb{R}$ с асимптотической дисперсией $\sigma^2(\theta)$. Пусть $\tau(\theta)$ дифференцируема на $\Theta \subset \mathbb{R}$.

Тогда $\tau(\hat{\theta}_n(X_1, \dots, X_n))$ - асимптотически нормальная оценка параметра $\tau(\theta)$ с асимптотической дисперсией $\sigma^2(\theta) \cdot (\tau'(\theta))^2$.

Доказательство.

По условию для любого $\theta \in \Theta$ выполнено

$$\sqrt{n} \left(\hat{\theta}_n(X_1, \dots, X_n) - \theta \right) \xrightarrow{d_\theta} \mathcal{N}(0, \sigma^2(\theta)).$$

Применяя дельта-метод, получаем, что

$$\sqrt{n} \left(\tau(\hat{\theta}_n(X_1, \dots, X_n)) - \tau(\theta) \right) \xrightarrow{d_\theta} \tau'(\theta) \cdot \mathcal{N}(0, \sigma^2(\theta)).$$

Остается заметить, что $\tau'(\theta) \cdot \mathcal{N}(0, \sigma^2(\theta)) = \mathcal{N}(0, \sigma^2(\theta) \cdot (\tau'(\theta))^2)$. ■

Замечание

Утверждения работают и в обратную сторону. Если мы умеем хорошо оценивать $\tau(\theta)$, существует обратная функция τ^{-1} и она непрерывна или дифференцируема.

Пример

Пусть X_1, \dots, X_n - выборка из экспоненциального распределения $\text{Exp}(\theta)$, $\theta > 0$. Найти асимптотически нормальную оценку θ .

Решение

Плотность X_i равна $p_\theta(x) = \theta \cdot e^{-\theta x} I\{x > 0\}$. Отсюда

$$E_\theta X_i = \frac{1}{\theta}, \quad D_\theta X_i = \frac{1}{\theta^2}.$$

Согласно ЦПТ

$$\sqrt{n} \left(\bar{X} - \frac{1}{\theta} \right) \xrightarrow{d_\theta} \mathcal{N} \left(0, \frac{1}{\theta^2} \right).$$

Положим $\tau(x) = \frac{1}{x}$. Тогда по утверждению 2

$$\sqrt{n} \left(\tau(\bar{X}) - \tau \left(\frac{1}{\theta} \right) \right) \xrightarrow{d_\theta} \tau' \left(\frac{1}{\theta} \right) \mathcal{N} \left(0, \frac{1}{\theta^2} \right)$$

или

$$\sqrt{n} \left(\frac{1}{\bar{X}} - \theta \right) \xrightarrow{d_\theta} \mathcal{N}(0, \theta^2).$$

Взаимоотношение свойств

Между свойствами есть следующие простые соотношения:

Сильная состоятельность \Rightarrow Состоятельность;

Асимптотическая нормальность \Rightarrow Состоятельность;

Методы построения оценок

Вопрос: какие можно предложить методы построения оценок с хорошими (прежде всего - асимптотическими) свойствами?

Принцип подстановки

Пусть для параметрического семейства $\{P_\theta, \theta \in \Theta\}$ нашелся такой функционал G , что для всех $\theta \in \Theta$ выполняется равенство

$$G(P_\theta) = \theta.$$

Тогда оценкой θ по методу подстановки называется

$$\theta^*(X_1, \dots, X_n) = G(P_n^*),$$

где P_n^* - эмпирическое распределение, построенное по выборке $\mathbf{X} = (X_1, \dots, X_n)$.

Вопрос: какие функционалы можно брать?

Идея: моменты определяют параметр распределения. То есть, у нас есть параметрическое семейство. Посчитаем моменты, которые являются функциями параметра. Если эти функции однозначно определяют параметр, то мы можем взять моменты в качестве функционала.

Метод моментов

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из неизвестного распределения $P \in \{P_\theta, \theta \in \Theta\}$, $\Theta \subset \mathbb{R}^k$. Пусть борелевские функции $g_1(x), \dots, g_k(x)$ таковы, что набор значений

$$E_\theta g_1(X_1), \dots, E_\theta g_k(X_1)$$

однозначно определяет параметр θ . Это означает, что функция $m : \Theta \rightarrow \mathbb{R}^k$, где

$$m_i(\theta) = E_\theta g_i(X_1), \quad i = 1, \dots, k,$$

является биекцией между Θ и $m(\Theta)$.

Определение 3.9. Оценкой θ по методу моментов с пробными функциями $g_1(x), \dots, g_k(x)$ называется $\theta^*(X_1, \dots, X_n)$ - решение следующей системы уравнений относительно θ :

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

Другими словами:

$$\theta^*(X_1, \dots, X_n) = m^{-1}\left(\overline{g_1(\mathbf{X})}, \dots, \overline{g_k(\mathbf{X})}\right).$$

Замечание (1)

Если множество Θ не является открытым, то вполне может оказаться, что вектор $\left(\overline{g_1(\mathbf{X})}, \dots, \overline{g_k(\mathbf{X})}\right)$ с большой положительной вероятностью не попадет в область определения функции m^{-1} . Так что разумно считать, что θ принимает все возможные значения.

Замечание (2)

Обычно используют стандартные пробные функции

$$g_1(x) = x, \dots, g_k(x) = x^k.$$

Пример

Пусть (X_1, \dots, X_n) - выборка из нормального распределения $\mathcal{N}(a, \sigma^2)$, $a \in \mathbb{R}$, $\sigma^2 > 0$, $\theta = (a, \sigma^2)$. Найти оценку θ по методу моментов.

Решение

Вспомним, что

$$E_{\theta}X_1 = a, \quad E_{\theta}X_1^2 = a^2 + \sigma^2.$$

Получаем систему

$$\begin{cases} a^* = \overline{X} \\ (a^*)^2 + (\sigma^*)^2 = \overline{X^2} \end{cases}$$

Решая, находим ответ:

$$a^* = \overline{X}, \quad (\sigma^*)^2 = \overline{X^2} - (\overline{X})^2 = S^2.$$

Теорема 3.1. (состоятельность оценки метода моментов)

Пусть $\theta_n^*(X_1, \dots, X_n) = m^{-1}\left(\overline{g_1(\mathbf{X})}, \dots, \overline{g_k(\mathbf{X})}\right)$ - корректно определенная оценка по методу моментов. Если функция m^{-1} непрерывна на множестве $m(\Theta)$, то θ_n^* - сильно состоятельная оценка параметра θ .

Доказательство.

Зафиксируем произвольное $\theta \in \Theta \subset \mathbb{R}^k$. Тогда согласно УЗБЧ для любого $j = 1, \dots, k$

$$\overline{g_j(\mathbf{X})} = \frac{1}{n} \sum_{i=1}^n g_j(X_i) \xrightarrow{P_{\theta-\text{п.н.}}} E_{\theta}g_j(X_1) = m_j(\theta).$$

По теореме о наследовании сходимости получаем:

$$\theta_n^*(X_1, \dots, X_n) = m^{-1}\left(\overline{g_1(\mathbf{X})}, \dots, \overline{g_k(\mathbf{X})}\right) \xrightarrow{P_{\theta-\text{п.н.}}} m^{-1}(m_1(\theta), \dots, m_k(\theta)) = m^{-1}(m(\theta)) = \theta$$

■

Упражнение

Если в дополнение к условиям теоремы для любого $\theta \in \Theta \subset \mathbb{R}^k$ и любого $j = 1, \dots, k$, конечно $E_{\theta}g_j^2(X_1)$, а функция m^{-1} дифференцируема на $m(\Theta)$, то для каждого $j = 1, \dots, k$ оценка $\theta_{j,n}^*(X_1, \dots, X_n)$ параметра θ_j по методу моментов является асимптотически нормальной.

Итак, мы увидели, что метод моментов поставяет состоятельные и асимптотически нормальные оценки параметра. Однако у него есть целый ряд недостатков:

- Подбор пробных функций. Если стандартные не подошли, то возникает проблема поиска подходящих функций. Например, рассмотрим распределение Коши со сдвигом: у нас есть выборка (X_1, \dots, X_n) , и мы взяли плотность распределения Коши со сдвигом $p_{\theta}(x) = \frac{1}{\pi(1 + (x - \theta)^2)}$, $\theta \in \mathbb{R}$. Тогда математическое ожидание $E_{\theta}X_1$ не определено, и, соответственно, ни одна из стандартных функций не может быть использована. Можно попробовать взять

$E_{\theta} \ln(|X_1| + 1)$. Такое математическое ожидание конечно, но метод моментов требует, чтобы мы могли посчитать эту функцию, как функцию от параметра. Это довольно сложно.

- Вторая трудность применения метода в том, что для нахождения оценки нужно знать функцию $m(\theta)$, что с практической точки зрения не представляется возможным.
- Получается, что применимость метода ограничена стандартными распределениями, для которых возможно аналитическое вычисление моментов, как функций от параметров.

На следующей лекции

Вернемся к примеру с распределением Коши со сдвигом. Забудем про идею с логарифмом. Что в таком случае можно предложить в качестве оценки параметра, если метод моментов не работает? Другая идея состоит в том, чтобы взять другой функционал. Для данного распределения θ не является математическим ожиданием, но оно является $\frac{1}{2}$ -квантилью, потому что плотность симметрична относительно θ , со-

ответственно в точке θ функция распределения будет достигать значения $\frac{1}{2}$. Можем

ли мы надеяться на то, что выборочная $\frac{1}{2}$ -квантиль хорошо приближает наш параметр θ ? На следующей лекции мы разберем теорему, которая доказывает, что выборочные квантили хорошо приближают настоящие, соответственно их можно использовать в качестве настоящих квантилей. Кроме того, в следующий раз мы посмотрим на принципы, по которым мы можем сравнивать оценки и разберем разные подходы к сравнению оценок.

4. Лекция 4

Сегодняшняя лекция будет посвящена продолжению изучения методов получения хороших оценок, а также принципам сравнению оценок.

Выборочные квантили

На прошлой лекции мы познакомились с методом моментов. Но есть распределения, где подбор нужных пробных функций крайне затруднен.

Например, пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из распределения Коши со сдвигом $\theta \in \mathbb{R}$: плотность

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

Все стандартные моменты у такого распределения либо не определены, либо бесконечны. Применение других функций выглядит безнадежным. Можно ли что-то предложить в качестве альтернативы методу моментов?

Идея: квантили определяют параметр распределения.

В примере θ является $\frac{1}{2}$ -квантилью распределения (медианой распределения).

Напомним, что если $\mathbf{X} = (X_1, \dots, X_n)$ - выборка, то выборочной p -квантилью называется p -квантиль эмпирического распределения P_n^* , то есть

$$Z_{n,p} = \begin{cases} X_{([np]+1)}, & \text{если } np \notin \mathbb{Z}; \\ X_{(np)}, & \text{если } np \in \mathbb{Z}. \end{cases}$$

Первой целью лекции будет доказательство теоремы о том, что выборочная квантиль является асимптотически нормальной оценкой настоящей квантили.

Теорема 4.1. (асимптотическая нормальность выборочной квантили)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка растущего размера из распределения с функцией распределения $F(x)$. Пусть $z_p, p \in (0, 1)$ - его p -квантиль, причем F дифференцируема в точке z_p и $F'(z_p) > 0$. Тогда утверждается, что $Z_{n,p}$ будет асимптотически нормальной оценкой настоящей квантили z_p , и выполнена следующая сходимость:

$$\sqrt{n}(Z_{n,p} - z_p) \xrightarrow{d} \mathcal{N}\left(0, \frac{p(1-p)}{(F'(z_p))^2}\right) \quad \text{при } n \rightarrow +\infty.$$

Доказательство.

Доказательство будет опираться на следующий факт.

Упражнение

Пусть дана последовательность биномиальных случайных величин $\xi_n \sim \text{Bin}(n, p_n)$, $n \in \mathbb{N}$, причем $p_n \rightarrow p \in (0, 1)$. Тогда

$$\frac{\xi_n - E\xi_n}{\sqrt{D\xi_n}} \xrightarrow{d} \mathcal{N}(0, 1).$$

Утверждение легко выводится из центральной предельной теоремы в схеме серий, но может быть доказано непосредственно методом характеристических функций.

Идея состоит в том, чтобы свести интересующую нас сходимости к данному упражнению. Для этого рассмотрим событие $\{X_{(k)} \leq x\}$ и заметим, что оно означает, что не менее k из случайных величин X_1, \dots, X_n попали в промежуток $(-\infty, x]$. Следовательно, его можно представить следующим образом:

$$\{X_{(k)} \leq x\} = \left\{ \sum_{i=1}^n I\{X_i \leq x\} \geq k \right\},$$

то есть хотя бы k элементов выборки попали на полуинтервал $(-\infty, x]$. Здесь уже видна связь необходимого нам события с упражнением, так как это сумма независимых индикаторов, то есть схема Бернулли, только в нашем случае вероятность того, что $X_i \leq x$ будет как-то зависеть от n . Запишем интересующую нас вероятность:

$$P(\sqrt{n}(X_{n,p} - z_p) \leq x) = P\left(Z_{n,p} \leq z_p + \frac{x}{\sqrt{n}}\right) = P\left(\sum_{i=1}^n I\left\{X_i \leq z_p + \frac{x}{\sqrt{n}}\right\} \geq np\right).$$

Обозначим

$$Y_n = \sum_{i=1}^n I\left\{X_i \leq z_p + \frac{x}{\sqrt{n}}\right\}.$$

Это биномиальная случайная величина со следующими характеристиками:

$$EY_n = nF\left(z_p + \frac{x}{\sqrt{n}}\right) \sim nF(z_p) = np,$$

$$DY_n = nF\left(z_p + \frac{x}{\sqrt{n}}\right)\left(1 - F\left(z_p + \frac{x}{\sqrt{n}}\right)\right) \sim np(1 - p).$$

Согласно упражнению получаем, что

$$\frac{Y_n - EY_n}{\sqrt{DY_n}} \xrightarrow{d} \mathcal{N}(0, 1).$$

Далее,

$$P(Y_n \geq np) = P\left(\frac{Y_n - EY_n}{\sqrt{DY_n}} \geq \frac{np - EY_n}{\sqrt{DY_n}}\right).$$

Найдем асимптотику выражения $\frac{np - EY_n}{\sqrt{DY_n}}$. Заметим, что

$$EY_n = n \cdot F\left(z_p + \frac{x}{\sqrt{n}}\right) = n\left(F(z_p) + \frac{x}{\sqrt{n}}F'(z_p) + o(n^{-\frac{1}{2}})\right) = np + xF'(z_p)\sqrt{n} + o(\sqrt{n}).$$

Значит,

$$\frac{np - EY_n}{\sqrt{DY_n}} = \frac{-xF'(z_p)\sqrt{n} + o(\sqrt{n})}{\sqrt{np(1-p)(1+o(1))}} = -\frac{xF'(z_p)}{\sqrt{p(1-p)}} + o(1).$$

Следовательно,

$$\lim_{n \rightarrow +\infty} P(\sqrt{n}(Z_{n,p} - z_p) \leq x) = 1 - \Phi\left(-\frac{xF'(z_p)}{\sqrt{p(1-p)}}\right) = \Phi\left(\frac{xF'(z_p)}{\sqrt{p(1-p)}}\right),$$

что и означает искомую сходимость:

$$\sqrt{n}(Z_{n,p} - z_p) \xrightarrow{d} \mathcal{N}\left(0, \frac{p(1-p)}{(F'(z_p))^2}\right) \quad \text{при } n \rightarrow +\infty.$$

■

Выборочная медиана

Определение 4.1. Медианой распределения называется его $\frac{1}{2}$ -квантиль. Выборочной медианой $\hat{\mu} = \hat{\mu}(X_1, \dots, X_n)$ выборки (X_1, \dots, X_n) называется величина

$$\hat{\mu}(X_1, \dots, X_n) = \begin{cases} X_{(k+1)}, & \text{если } n = 2k + 1; \\ \frac{X_{(k)} + X_{(k+1)}}{2}, & \text{если } n = 2k. \end{cases}$$

Следствие

В условиях теоремы об асимптотической нормальности выборочной квантили выполнено

$$\sqrt{n}(\hat{\mu} - z_{1/2}) \xrightarrow{d} \mathcal{N}\left(0, \frac{1}{4(F'(z_{1/2}))^2}\right).$$

Доказательство.

Отметим, что доказательство теоремы остается верным, если вместо $Z_{n,p}$ рассматривать любую порядковую статистику $X_{(k)}$ при $k = np + O(1)$, то есть

$$\sqrt{n}(X_{(k)} - z_p) \xrightarrow{d} \mathcal{N}\left(0, \frac{p(1-p)}{(F'(z_p))^2}\right) \quad \text{при } n \rightarrow +\infty.$$

Отметим, что для любого $x \in \mathbb{R}$ по подпоследовательности четных чисел $n = 2k$ выполнено

$$\begin{aligned} \lim_{k \rightarrow +\infty} P(\sqrt{n}(X_{(k+1)} - z_{1/2}) \leq x) &\leq \\ &\leq \lim_{k \rightarrow +\infty} P(\sqrt{n}(\hat{\mu} - z_{1/2}) \leq x) \leq \\ &\leq \lim_{k \rightarrow +\infty} P(\sqrt{n}(X_{(k)} - z_{1/2}) \leq x), \end{aligned}$$

что и доказывает искомую сходимость, так как крайние пределы равны. ■

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из распределения Коши со сдвигом $\theta \in \mathbb{R}$: плотность

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

Заметим, что наше распределение подходит под условия теоремы: плотность непрерывна и положительна в точке θ . Тогда θ - это медиана распределения и, стало быть, по теореме выборочная медиана будет асимптотически нормальной (и состоятельной!) оценкой параметра:

$$\sqrt{n}(\hat{\mu} - \theta) \xrightarrow{d} \mathcal{N}\left(0, \frac{1}{4p_\theta^2(\theta)}\right) = \mathcal{N}\left(0, \frac{\pi^2}{4}\right).$$

Сравнение оценок

Мы уже предъявили несколько методов, которые позволяют находить оценки с хорошими свойствами. Пусть $\hat{\theta}(\mathbf{X})$, $\theta^*(\mathbf{X})$ - две хорошие оценки параметра $\theta \in \Theta$. Как понять, какая из них лучше оценивает θ ? Для этого нужны способы сравнения оценок. В рамках курса мы разберем 4 подхода и будем решать задачи, связанные с нахождением наилучших оценок в этих подходах. 3 из них основаны на работе с функцией риска, 4ый подход связан со сравнением асимптотически нормальных оценок. В первую очередь поймем, какие могут быть способы определения близости оценки и параметра.

Определение 4.2. Борелевская неотрицательная функция $\rho(x, y) \geq 0$ называется функцией потерь. Если $\theta^*(\mathbf{X})$, то ее величиной потерь называется $\rho(\theta^*(\mathbf{X}), \theta)$.

Примеры

- $\rho(x, y) = |x - y|$;

- $\rho(x, y) = (x - y)^2$ - квадратичная функция потерь;
- если $\theta \in \mathbb{R}^d$, $d > 1$, то для неотрицательно определенной матрицы $A \in \text{Mat}(d \times d)$ можно определить

$$\rho(x, y) = \langle A(x - y), x - y \rangle$$

Функция риска

Сама по себе функция потерь мало что нам дает, так как она зависит от параметра, зависит от x и является случайной. Поэтому введем понятие функции риска, которая позволяет избавиться хотя-бы от одной зависимости.

Определение 4.3. При заданной функции потерь $\rho(x, y) \geq$ функцией риска оценки $\theta^*(\mathbf{X})$ называется

$$R_{\theta^*}(\theta) = E_{\theta} \rho(\theta^*(\mathbf{X}), \theta).$$

Базовое сравнение оценок, в первую очередь, основано на сравнении их функций риска. Здесь есть сразу несколько подходов, мы рассмотрим три основных: равномерный, байесовский и минимаксный.

Равномерный подход

Правило сравнения. Считаем, что оценка $\theta^*(\mathbf{X})$ лучше оценивает параметр θ , чем оценка $\hat{\theta}(\mathbf{X})$, если для любого $\theta \in \Theta$ выполнено

$$R_{\theta^*}(\theta) \leq R_{\hat{\theta}}(\theta),$$

причем для некоторого $\theta \in \Theta$ неравенство строгое.

Определение 4.4. Оценка $\theta^*(\mathbf{X})$ называется наилучшей в классе \mathcal{K} , если она лучше любой другой оценки $\hat{\theta}(\mathbf{x}) \in \mathcal{K}$.

Почему необходимо рассматривать именно классы? Оказывается, что постановка задачи сравнения во всем классе оценок без ограничений бессмысленна.

Замечание

Наилучшая оценка не всегда существует. Например, если \mathcal{K} - это класс всех возможных оценок, $\theta \in \mathbb{R}$ и $\rho(x, y) = (x - y)^2$ - это квадратичная функция потерь, то наилучшей оценки нет.

Действительно, можно для любого $\theta_0 \in \Theta$ рассмотреть оценку $\hat{\theta}_0(\mathbf{X}) \equiv \theta_0$, то ее функция риска в точке θ_0 обращается в ноль:

$$R_{\hat{\theta}_0}(\theta_0) = E_{\theta_0} \left(\hat{\theta}_0(\mathbf{X}) - \theta_0 \right)^2 = 0.$$

Тем самым, функция риска потенциальной наилучшей оценки равна тождественно нулю, то есть сама она равна всем константам одновременно. Противоречие.

Отсюда возникает вопрос: какие классы разумно рассматривать? Задача о поиске наилучшей оценки в равномерном подходе имеет смысл в классе оценок с одинаковым математическим ожиданием. именно в нем мы будем ее решать.

Определение 4.5. Равномерный подход с квадратичной функцией потерь называется среднеквадратичным.

Определение 4.6. Оценка $\theta^*(\mathbf{X})$ называется допустимой оценкой θ , если нет другой оценки $\hat{\theta}_0$, которая оценивает θ лучше, чем $\theta^*(\mathbf{X})$, в равномерном подходе.

Байесовский подход

Правило сравнения. Пусть Q - некоторое заданное распределение вероятностей на Θ (априорное распределение). Тогда для оценки $\theta^*(\mathbf{X})$ с функцией риска $R_{\theta^*}(\theta)$ определим

$$\rho_Q(\theta^*) = E_Q R_{\theta^*}(\theta) = \int_{\Theta} R_{\theta^*}(t) Q(dt).$$

Например, если Q имеет плотность $q(t)$ по мере λ , то

$$\rho_Q(\theta^*) = \int_{\Theta} R_{\theta^*}(t) q(t) \lambda(dt).$$

Определение 4.7. Оценка $\theta^*(\mathbf{X})$ называется наилучшей оценкой θ в байесовском подходе для заданного априорного распределения Q (байесовской оценкой), если

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

Байесовский подход интересен тем, что мы получаем серию оценок. Понятно, что наилучшая оценка зависит от априорного распределения Q . Эти оценки интересны тем, что они являются допустимыми.

Упражнение

Пусть оценка $\theta^*(\mathbf{X})$ является наилучшей оценкой θ в байесовском подходе для некоторого априорного распределения Q . Тогда $\theta^*(\mathbf{X})$ - это допустимая оценка θ .

Минимаксный подход

Правило сравнения. Пусть $\theta^*(\mathbf{X})$ - оценка θ с функцией риска $R_{\theta^*}(\theta)$. Тогда определим

$$\rho_{\max}(\theta^*) = \sup_{\theta \in \Theta} R_{\theta^*}(\theta).$$

Определение 4.8. Оценка $\theta^*(\mathbf{X})$ называется наилучшей оценкой θ в минимаксном подходе (минимаксной оценкой), если

$$\rho_{\max}(\theta^*) = \inf_{\hat{\theta}(\mathbf{X})} \rho_{\max}(\hat{\theta}),$$

то есть у $\theta^*(\mathbf{X})$ наименьшее максимальное значение функции риска.

Как мы увидим в дальнейшем, минимаксный и байесовский подходы тесно связаны.

Проиллюстрируем на чем основаны подходы. Предположим, что $\Theta = [0, 1]$, априорное распределение $Q = \cup[0, 1]$ (равномерное распределение на отрезке), и нарисуем примеры двух функций риска. Попытаемся сравнить оценки θ^* и $\hat{\theta}$ во всех подходах.

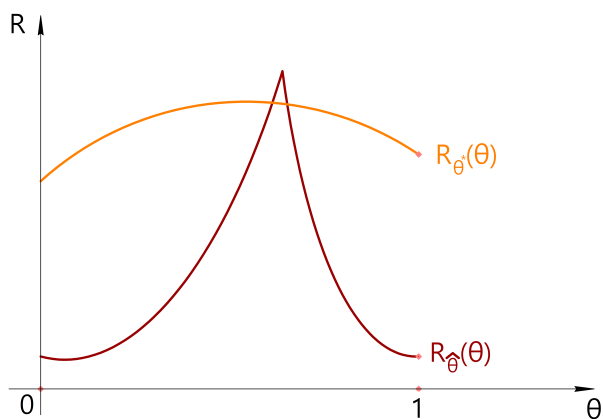


Рис. 3

Мы можем сделать следующие выводы:

- θ^* лучше в минимаксном подходе, так как у нее максимальное значение функции риска ниже;
- $\hat{\theta}$ в свою очередь лучше в байесовском подходе, так как интегральное значение — это площадь графика;
- В равномерном подходе оценки не сравнимы.

Асимптотический подход

Это принципиально другой подход к сравнению оценок. Он апеллирует к асимптотике, то есть мы имеем дело с последовательностями асимптотически нормальных оценок. Этот подход практически никак не связан с предыдущими и не использует функции риска.

Правило сравнения. Пусть $\hat{\theta}_{n,1}(X_1, \dots, X_n)$ и $\hat{\theta}_{n,2}(X_1, \dots, X_n)$ - две асимптотически нормальные оценки $\theta \in \mathbb{R}$ с асимптотическими дисперсиями $\sigma_1^2(\theta)$ и $\sigma_2^2(\theta)$. Тогда оценка $\hat{\theta}_{n,1}(X_1, \dots, X_n)$ лучше оценивает параметр θ , чем оценка $\hat{\theta}_{n,2}(X_1, \dots, X_n)$, если для любого $\theta \in \Theta$ выполнено

$$\sigma_1^2(\theta) \leq \sigma_2^2(\theta),$$

причем для некоторого $\theta \in \Theta$ неравенство строгое.

Определение 4.9. Оценка $\hat{\theta}_n(X_1, \dots, X_n)$ называется наилучшей в асимптотическом подходе, если она лучше любой другой оценки.

Замечание

Мы увидим, что получить совсем честное решение здесь невозможно. Какова бы ни была асимптотически нормальная оценка с положительной асимптотической дисперсией, всегда можно подобрать другую оценку, у которой асимптотическая дисперсия будет меньше хотя бы в одной точке.

Пример

Пусть X_1, \dots, X_n - выборка из $\mathcal{N}(\theta, 1)$, $\theta \in \mathbb{R}$. Сравните в асимптотическом подходе оценки \bar{X} и $\hat{\mu}$.

Решение

Обе оценки являются асимптотически нормальными. Асимптотическая нормальность \bar{X} следует из ЦПТ

$$\sqrt{n}(\bar{X} - \theta) \xrightarrow{d} \mathcal{N}(0, 1).$$

По теореме об асимптотической нормальности выборочной квантили получаем:

$$\sqrt{n}(\hat{\mu} - \theta) \xrightarrow{d} \mathcal{N}\left(0, \frac{1}{4p_\theta^2(\theta)}\right) = \mathcal{N}\left(0, \frac{\pi}{2}\right).$$

Вывод: выборочное среднее лучше.

Дальнейшая цель: научиться строить наилучшие оценки в разных подходах.

5. Лекция 5

Сегодняшняя лекция будет посвящена неравенству Рао-Крамера. Мы познакомимся с тем, что такое информация Фишера и эффективные оценки, и получим одно из частичных решений нахождения наилучшей оценки в равномерном подходе.

Доминируемые семейства

Далее в курсе мы часто будем предполагать, что имеем дело с так называемыми доминируемыми параметрическими семействами, имеющими (обобщенные) плотности.

Определение 5.1. Пусть $\{P_\theta, \theta \in \Theta\}$ - параметрическое семейство распределений на выборочном пространстве $(\mathcal{X}, \mathcal{B}_\mathcal{X})$. Если для любого $\theta \in \Theta$ мера P_θ имеет плотность $p_\theta(x)$ по одной и той же (может быть, σ -конечной) мере μ на $(\mathcal{X}, \mathcal{B}_\mathcal{X})$, то семейство $\{P_\theta, \theta \in \Theta\}$ называется доминируемым относительно μ .

Напомним, что плотность понимается в смысле производной Радона-Никодима:

$$p_\theta(\theta) = \frac{dP_\theta}{d\mu}, \quad P_\theta(B) = \int_B p_\theta(x) \mu(dx).$$

Примеры:

- 1) Если μ - это мера Лебега на \mathbb{R} , то P_θ - это абсолютно непрерывные распределения. В этом случае обобщенная плотность - это обычная плотность из курса теории вероятностей. Например, в семействе нормальных распределений $\mathcal{N}(\theta, 1)$, $\theta \in \mathbb{R}$, плотность равна

$$p_\theta(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\theta)^2}{2}}.$$

- 2) Если μ - это считающая мера (мера, которая сопоставляет каждому множеству количество целых точек внутри него) на \mathbb{Z} , то P_θ - это дискретные распределения на \mathbb{Z} . В этом случае обобщенная плотность - это вероятность получить данное значение. Например, в семействе пуассоновских распределений $\text{Pois}(\theta)$, $\theta > 0$, плотность равна

$$p_\theta(x) = P_\theta(\{x\}) = \frac{\theta^x}{x!} e^{-\theta} \mathbf{I}\{x \in \mathbb{Z}_+\}.$$

Условия регулярности

Мы разберем частичное решение в среднеквадратичном подходе (равномерный подход с квадратичной функцией потерь).

Пусть \mathbf{X} - наблюдение с неизвестным распределением $P \in \{P_\theta, \theta \in \Theta\}$, где $\{P_\theta, \theta \in \Theta\}$ - доминируемое семейство относительно меры μ с плотностью $p_\theta(x)$. Введем следующие условия регулярности:

- 1) $\Theta \subset \mathbb{R}$ - открытый интервал (может быть, бесконечный).
- 2) Множество $A = \{x : p_\theta(x) > 0\}$ не зависит от θ (носитель плотности не зависит от параметра). Это условие уже отбрасывает некоторые семейства распределений, например, под него не попадает семейство равномерное распределение на отрезке $[0, \theta]$, потому что носитель зависит от параметра.
- 3) Для любой статистики $S(\mathbf{X})$ с равномерно ограниченным вторым моментом, $E_\theta S^2(\mathbf{X}) < M < +\infty$ для всех $\theta \in \Theta$, выполнено: для любого $\theta \in \Theta$

$$\frac{\partial}{\partial \theta} E_\theta S(\mathbf{X}) = E_\theta \left(S(\mathbf{X}) \cdot \frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) \right).$$

- 4) Величина

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

положительна и конечна для всех $\theta \in \Theta$.

Обсудим подробнее условие 3:

Оно означает, что мы можем дифференцировать под знаком интеграла. Действительно

$$\begin{aligned} \frac{\partial}{\partial \theta} E_\theta S(\mathbf{X}) &= \frac{\partial}{\partial \theta} \int_A S(\mathbf{x}) p_\theta(\mathbf{x}) \mu(d\mathbf{x}) = \int_A S(\mathbf{x}) \frac{\partial}{\partial \theta} p_\theta(\mathbf{x}) \mu(d\mathbf{x}) = \\ &= \int_A S(\mathbf{x}) \frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{x}) \cdot p_\theta(\mathbf{x}) \mu(d\mathbf{x}) = E_\theta \left(S(\mathbf{X}) \cdot \frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) \right) \end{aligned}$$

Вклад и информация Фишера

Величины, участвующие в условии 4, имеют специальные названия.

Определение 5.2. Случайная величина

$$U_{\theta}(\mathbf{X}) = \frac{\partial}{\partial \theta} \ln p_{\theta}(\mathbf{X})$$

называется вкладом наблюдения \mathbf{X} .

Определение 5.3. Функция $I_{\mathbf{X}}(\theta)$ называется количеством информации о параметре θ , содержащимся в наблюдении \mathbf{X} (информации по Фишеру).

$$I_{\mathbf{X}}(\theta) = E_{\theta}(U_{\theta}(\mathbf{X}))^2.$$

Неравенство Рао-Крамера

Оказывается, что в случае регулярных семейств есть нижняя граница для дисперсии у несмещенной оценки некоторой функции от параметра. Эта граница постулируется утверждением, которое называется неравенством Рао-Крамера

Теорема 5.1. (неравенство Рао-Крамера)

Пусть выполнены условия регулярности 1-4. Если $\hat{\theta}(\mathbf{X})$ - это несмещенная оценка $\tau(\theta)$ с равномерно ограниченным вторым моментом $E_{\theta}(\hat{\theta}(\mathbf{X}))^2 < M < +\infty$ для любого $\theta \in \Theta$, то для всех $\theta \in \Theta$ выполняется неравенство

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

Доказательство.

Воспользуемся третьим условием регулярности, подставив в нее $S(\mathbf{X}) \equiv 1$:

$$\frac{\partial}{\partial \theta} E_{\theta} S(\mathbf{X}) = \frac{\partial}{\partial \theta} E_{\theta} 1 = 0, \quad E_{\theta}(S(\mathbf{X})U_{\theta}(\mathbf{X})) = E_{\theta}U_{\theta}(\mathbf{X}).$$

Следовательно, $E_{\theta}U_{\theta}(\mathbf{X}) = 0$.

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

$$\begin{aligned} \frac{\partial}{\partial \theta} E_{\theta} S(\mathbf{X}) &= \frac{\partial}{\partial \theta} E_{\theta} \hat{\theta}(\mathbf{X}) = \tau'(\theta), \\ E_{\theta}(S(\mathbf{X})U_{\theta}(\mathbf{X})) &= E_{\theta}(\hat{\theta}(\mathbf{X})U_{\theta}(\mathbf{X})) = E_{\theta}((\hat{\theta}(\mathbf{X}) - \tau(\theta))U_{\theta}(\mathbf{X})). \end{aligned}$$

Тем самым, $\tau'(\theta) = E_{\theta}((\hat{\theta}(\mathbf{X}) - \tau(\theta))U_{\theta}(\mathbf{X}))$. Применяя неравенство Коши-Буняковского, получаем искомое соотношение:

$$(\tau'(\theta))^2 \leq E_{\theta}(\hat{\theta}(\mathbf{X}) - \tau(\theta))^2 \cdot E_{\theta}U_{\theta}^2(\mathbf{X}) = D_{\theta}\hat{\theta}(\mathbf{X}) \cdot I_{\mathbf{X}}(\theta).$$



Следствие 1

Если в условиях неравенства Рао-Крамера выполнено, что $\tau(\theta) = \theta$, то для всех $\theta \in \Theta$

$$D_{\theta}\hat{\theta} \geq \frac{1}{I_{\mathbf{X}}(\theta)}.$$

Вопрос: чему равняется сама информация?

Утверждение

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - это выборка. Обозначим через

$$i(\theta) = E_{\theta} \left(\frac{\partial}{\partial \theta} \ln p_{\theta}(X_1) \right)^2$$

информацию Фишера одного наблюдения. Тогда

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

Следствие 2

Если в условиях неравенство Рао-Крамера $\mathbf{X} = (X_1, \dots, X_n)$ - это выборка и $\tau(\theta)$ не зависит от n , то

$$D_{\theta}\hat{\theta}(\mathbf{X}) \geq \frac{(\tau'(\theta))^2}{n \cdot i(\theta)} = \Omega\left(\frac{1}{n}\right).$$

Вопрос: возможно ли равенство в неравенстве Рао-Крамера?

Определение 5.4. Пусть $\hat{\theta}(\mathbf{X})$ - несмещенная оценка $\tau(\theta)$. Будем называть ее эффективной оценкой $\tau(\theta)$, если для нее достигается равенство в неравенстве Рао-Крамера, то есть для всех $\theta \in \Theta$

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

Критерий эффективности

Лемма (критерий эффективности)

В условиях неравенства Рао-Крамера $\hat{\theta}(\mathbf{X})$ будет эффективной оценкой $\tau(\theta)$ тогда и только тогда, когда для любого $\theta \in \Theta$ Р $_{\theta}$ -п.н. выполняется равенство

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

Доказательство.

Неравенство Рао-Крамера получалось применением неравенства Коши-Буняковского к равенству:

$$E_{\theta} \left((\hat{\theta}(\mathbf{X}) - \tau(\theta)) U_{\theta}(\mathbf{X}) \right) = \tau'(\theta).$$

Как известно, в неравенстве Коши-Буняковского равенство достигается тогда и только тогда, когда между случайными величинами есть линейная зависимость с вероятностью 1. Тут надо понимать, что поскольку мы фиксируем θ и применяем неравенство, то линейная зависимость между $\hat{\theta}(\mathbf{X}) - \tau(\theta)$ и $U_{\theta}(\mathbf{X})$ для каждого θ своя. Следовательно, $\hat{\theta}(\mathbf{X})$ будет эффективной оценкой для $\tau(\theta)$ только в том случае, если существует такая функция $c(\theta)$, что

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

Найдем $c(\theta)$. Домножим обе части последнего равенства на $U_{\theta}(\mathbf{X})$ и возьмем математическое ожидание:

$$E_{\theta} \left((\hat{\theta}(\mathbf{X}) - \tau(\theta)) U_{\theta}(\mathbf{X}) \right) = E_{\theta} (c(\theta) U_{\theta}^2(\mathbf{X})).$$

Но левая часть есть $\tau'(\theta)$, а второй момент для вклада - информация Фишера. Тогда

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

■

Разберем пример, который показывает, как необходимо подходить к задачам такого рода. Идея состоит в том, что необходимо вычислить вклад и посмотреть, представляется ли он в виде произведения функции от параметра на разность функции от x и функции от параметра. Если такое представление возможно, то функция от x будет эффективной оценкой того, что из нее вычитается, а информацию Фишера можно посчитать с помощью равенства, в котором наш коэффициент выражается через производную оцениваемой функции и информацию Фишера.

Пример

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 $i(\theta)$ одного наблюдения.

Решение

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

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

Теперь несложно найти вклад выборки:

$$\ln p_{\theta}(\mathbf{X}) = \sum_{i=1}^n (X_i \ln \theta + (1 - X_i) \ln(1 - \theta)),$$

$$U_{\theta}(\mathbf{X}) = \frac{X_1 + \dots + X_n}{\theta} - \frac{n - X_1 - \dots - X_n}{1 - \theta} = n \left(\frac{\bar{X}}{\theta} - \frac{1 - \bar{X}}{1 - \theta} \right) = \frac{n}{\theta(1 - \theta)} (\bar{X} - \theta).$$

Отсюда мы можем сделать вывод, что \bar{X} - эффективная оценка для $\tau(\theta) = \theta$. Значит $\tau'(\theta) = 1$ и

$$\frac{I_{\mathbf{X}}(\theta)}{\tau'(\theta)} = c^{-1}(\theta) = \frac{n}{\theta(1 - \theta)}$$

Отсюда получаем информацию Фишера:

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

Выводы

- 1) Эффективная оценка для заданного семейства существует максимум для одной функции $\tau(\theta)$ (с точностью до умножения и добавления констант).
- 2) Требуется регулярность семейства распределений.
- 3) Если эффективная оценка существует, то плотность \mathbf{X} имеет специальный вид:

$$p_{\theta}(\mathbf{X}) = h(\mathbf{X}) e^{A(\theta)\hat{\theta}(\mathbf{X}) + B(\theta)},$$

где функции A, B не зависят от \mathbf{X} , а h - не зависит от θ .

- 4) Эффективная оценка $\tau(\theta)$ - это несмещенная оценка $\tau(\theta)$ с равномерно наименьшей дисперсией. Значит она - наилучшая оценка θ в среднеквадратичном подходе в классе всех несмещенных оценок $\tau(\theta)$:

$$E_{\theta}(\hat{\theta}(\mathbf{X}) - \theta)^2 = D_{\theta}\hat{\theta}(\mathbf{X}) + (\tau(\theta) - \theta)^2.$$

О выполнении условий регулярности

Естественно задаться вопросом о выполнении условий регулярности. Условия 1, 2 и 4 не вызывают вопросов, но для доказательства неравенства Рао-Крамера принципиально условие 3, прямая проверка которого очевидно затруднительная. Однако имеет место следующая теорема.

Теорема 5.2. Пусть для почти всех x по мере μ функция $\sqrt{p_\theta(x)}$ непрерывно дифференцируема по θ , а также информация Фишера $I_{\mathbf{X}}(\theta)$ конечна, положительна и непрерывна по θ . Тогда выполнено условие регулярности 3.

Доказательство теоремы можно найти в литературе, мы его опустим и приведем несколько замечаний.

Замечание (1)

Для конкретной оценки $\hat{\theta}(\mathbf{X})$ для выполнения неравенства Рао-Крамера достаточно конечности второго момента и возможности дифференцирования по параметру под знаком интеграла. Последнее можно проверить отдельно.

Замечание (2)

Условие равномерной ограниченности второго момента не является ограничительным. Если в условиях регулярности оценка $\hat{\theta}(\mathbf{X})$ такова, что функция $E_\theta(\hat{\theta}(\mathbf{X}))^2$ непрерывна по θ , то она будет ограничена на любом отрезке внутри Θ . Сузив значение параметра до интервала внутри отрезка, мы получим неравенство Рао-Крамера на интервале. Но раз это верно для любого интервала внутри Θ , то оно будет выполнено и на всем Θ .

Информация Фишера статистик

Понятие информации Фишера можно ввести не только для самого наблюдения, но и для статистик от него, чье семейство распределений является доминируемым.

Определение 5.5. Пусть $S(\mathbf{X})$ - это некоторая статистика, имеющая плотность $g_\theta(s)$ по мере λ , то есть семейство ее распределений $(G_\theta, \theta \in \Theta)$ является доминируемым относительно λ . Тогда информацией Фишера статистики $S(\mathbf{X})$ называется

$$I_S(\theta) = E_\theta \left(\left(\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) \right)^2 \right).$$

У информации Фишера есть несколько свойств, которые частично объясняют ее название.

Свойство (1)

Если распределение $S(\mathbf{X})$ не зависит от θ , то $I_S(\theta) = 0$.

Доказательство.

очевидно, так как в таком случае $\ln g_\theta(S(\mathbf{X}))$ не будет зависеть от θ и при дифференцировании по нему обратится в 0. ■

Свойство (2)

Пусть $S(\mathbf{X})$ и $T(\mathbf{X})$ - две независимые статистики для всех $\theta \in \Theta$. Обозначим $H(\mathbf{X}) = (S(\mathbf{X}), T(\mathbf{X}))$. Тогда в условиях регулярности выполнено равенство

$$I_h(\theta) = I_S(\theta) + I_T(\theta).$$

Доказательство.

Пусть статистика $S(\mathbf{X})$ имеет плотность $g_\theta(s)$ по мере λ_1 , а статистика $T(\mathbf{X})$ - плотность $f_\theta(t)$ по мере λ_2 . Тогда случайный вектор $H(\mathbf{X}) = (S(\mathbf{X}), T(\mathbf{X}))$ имеет совместную плотность $h_\theta(s, t) = g_\theta(s)f_\theta(t)$ по мере $\lambda_1 \times \lambda_2$. Тогда

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

Следовательно, в силу условий регулярности

$$E_\theta \left(\frac{\partial}{\partial \theta} \ln h_\theta(H(\mathbf{X})) \right) = E_\theta \left(\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) \right) + E_\theta \left(\frac{\partial}{\partial \theta} \ln f_\theta(T(\mathbf{X})) \right) = 0.$$

Но тогда в силу независимости

$$\begin{aligned} I_H(\theta) &= D_\theta \left(\frac{\partial}{\partial \theta} \ln h_\theta(H(\mathbf{X})) \right) = D_\theta \left(\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) + \frac{\partial}{\partial \theta} \ln f_\theta(T(\mathbf{X})) \right) = \\ &= D_\theta \left(\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) \right) + D_\theta \left(\frac{\partial}{\partial \theta} \ln f_\theta(T(\mathbf{X})) \right) = I_S(\theta) + I_T(\theta). \end{aligned}$$

■

Свойство (3)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - это выборка. Обозначим через

$$i(\theta) = E_\theta \left(\frac{\partial}{\partial \theta} \ln p_\theta(X_1) \right)^2$$

информацию Фишера одного наблюдения. Тогда

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

Доказательство.

Доказательство следует из свойства (2) по индукции.

■

Свойство (4)

В условиях регулярности для любой статистики $S(\mathbf{X})$ выполнено: для любого $\theta \in \Theta$

$$I_S(\theta) \leq I_{\mathbf{X}}(\theta).$$

Доказательство.

Доказательство этого свойства использует свойство условного математического ожидания (будет рассмотрено позже). Рекомендуются вернуться к доказательству этого свойства после изучения свойства математического ожидания.

Пусть $p_\theta(x)$ - это плотность \mathbf{X} по мере μ , а $g_\theta(s)$ - плотность статистики $S(\mathbf{X})$ по мере λ . Покажем, что

$$E_\theta \left(\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) \middle| S(\mathbf{X}) \right) = \frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})).$$

Левая часть выражения означает, что мы берем условное математическое ожидание от $\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X})$ по статистике $S(\mathbf{X})$.

Проверим по определению. Правая часть должна удовлетворять двум свойствам: интегральному и быть $S(\mathbf{X})$ -измеримой. Правая часть $S(\mathbf{X})$ -измерима, поэтому остается проверить интегральное свойство. Для любого борелевского множества B надо проверить, что

$$E_\theta \left(\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) \cdot I\{S(\mathbf{X}) \in B\} \right) = E_\theta \left(\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) \cdot I\{S(\mathbf{X}) \in B\} \right).$$

Посчитаем левую часть. Обозначим $A = \{x : p_\theta(x) > 0\}$. Тогда

$$\begin{aligned} E_\theta \left(\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) \cdot I\{S(\mathbf{X}) \in B\} \right) &= \int_A \frac{\partial \ln p_\theta(\mathbf{x})}{\partial \theta} p_\theta(\mathbf{x}) I\{S(\mathbf{x}) \in B\} \mu(d\mathbf{x}) = \\ &= \int_A \frac{\partial p_\theta(\mathbf{x})}{\partial \theta} I\{S(\mathbf{x}) \in B\} \mu(d\mathbf{x}) = \frac{\partial}{\partial \theta} \int_A p_\theta(\mathbf{x}) I\{S(\mathbf{x}) \in B\} \mu(d\mathbf{x}) = \frac{\partial}{\partial \theta} P_\theta(S(\mathbf{X}) \in B). \end{aligned}$$

Но с другой стороны

$$\begin{aligned} \frac{\partial}{\partial \theta} P_\theta(S(\mathbf{X}) \in B) &= \frac{\partial}{\partial \theta} \int_B g_\theta(s) \lambda(ds) = \int_B \frac{\partial \ln g_\theta(s)}{\partial \theta} g_\theta(s) \lambda(ds) = \\ &= E_\theta \left(\frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})) \cdot I\{S(\mathbf{X}) \in B\} \right). \end{aligned}$$

Мы получили представление для условного математического ожидания. Теперь необходимо доказать неравенство для информации. Рассмотрим величину

$$M(\theta) = E_\theta \left(\frac{\partial \ln p_\theta(\mathbf{X})}{\partial \theta} \cdot \frac{\partial \ln g_\theta(S(\mathbf{X}))}{\partial \theta} \right).$$

Согласно неравенству Коши-Буняковского

$$M^2(\theta) \leq E_{\theta} \left(\left(\frac{\partial}{\partial \theta} \ln p_{\theta}(\mathbf{X}) \right)^2 \right) E_{\theta} \left(\left(\frac{\partial}{\partial \theta} \ln g_{\theta}(S(\mathbf{X})) \right)^2 \right) = I_{\mathbf{X}}(\theta) \cdot I_S(\theta).$$

Теперь докажем, что $M(\theta) = I_s(\theta)$. Для ого заметим, что по формуле полной вероятности

$$E_{\theta} \left(\frac{\partial \ln p_{\theta}(\mathbf{X})}{\partial \theta} \cdot \frac{\partial \ln g_{\theta}(S(\mathbf{X}))}{\partial \theta} \right) = E_{\theta} \left(E_{\theta} \left(\frac{\partial \ln p_{\theta}(\mathbf{X})}{\partial \theta} \cdot \frac{\partial \ln g_{\theta}(S(\mathbf{X}))}{\partial \theta} \right) \middle| S(\mathbf{X}) \right).$$

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

$$\begin{aligned} E_{\theta} \left(\frac{\partial \ln p_{\theta}(\mathbf{X})}{\partial \theta} \cdot \frac{\partial \ln g_{\theta}(S(\mathbf{X}))}{\partial \theta} \right) &= E_{\theta} \left(\frac{\partial \ln g_{\theta}(S(\mathbf{X}))}{\partial \theta} E \left(\frac{\partial \ln p_{\theta}(\mathbf{X})}{\partial \theta} \middle| S(\mathbf{X}) \right) \right) = \\ &= E_{\theta} \left(\left(\frac{\partial \ln g_{\theta}(S(\mathbf{X}))}{\partial \theta} \right)^2 \right) = I_S(\theta). \end{aligned}$$

Тогда

$$I_S^2(\theta) \leq I_{\mathbf{X}}(\theta) I_S(\theta),$$

что и дает искомое неравенство

$$I_S(\theta) \leq I_{\mathbf{X}}(\theta).$$

■

Свойство (5)

В условиях регулярности $I_s(\theta) = I_{\mathbf{X}}(\theta)$ для любого $\theta \in \Theta$ тогда и только тогда, когда статистика $S(\mathbf{X})$ является достаточной для семейства распределений $(P, \theta \in \Theta)$.

Данное свойство мы докажем позднее, когда пройдем достаточные статистики.

Приведем пример того, как можно вычислять информации Фишера статистик от нашей выборки.

Пример

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

Решение

Начнем с того, что найдем распределение $X_{(1)}$. Для этого заметим, что для $x > 0$

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

Следовательно, $X_{(1)} \sim \text{Exp}(n\theta)$. Теперь найдем информацию Фишера. Начнем с $I_{\mathbf{X}}(\theta)$. Вспомним, что $I_{\mathbf{X}}(\theta) = n \cdot i(\theta)$. 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

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

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

$$I_{X_{(1)}} = E_{\theta} \left(nX_{(1)} - \frac{1}{\theta} \right)^2 = n^2 D_{\theta} X_{(1)} = \frac{1}{\theta^2}.$$

6.

Лекция 6

Многомерный вариант

На прошлой лекции мы получили нижнюю оценку дисперсии несмещенной оценки функции одномерного параметра. Пусть теперь \mathbf{X} - наблюдение с неизвестным распределением из доминируемого параметрического семейства $\{P_{\theta}, \theta \in \Theta\}$ с плотностью $p_{\theta}(x)$, где $\theta \in \Theta \subset \mathbb{R}^k$, $k > 1$ - многомерный параметр. Как можно определить информацию Фишера в подобной ситуации?

Введем следующие обозначения:

- 1) $\theta = (\theta_1, \dots, \theta_k)^T$ (вектор-столбец);
- 2) вклад наблюдения \mathbf{X} - вектор-столбец частных производных логарифма плотности:

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

- 3) информационная матрица размера $k \times k$:

$$I_{\mathbf{X}}(\theta) = \left(E_{\theta} \left(\frac{\partial \ln p_{\theta}(\mathbf{X})}{\partial \theta_i} \cdot \frac{\partial \ln p_{\theta}(\mathbf{X})}{\partial \theta_j} \right), i, j = 1, \dots, k \right).$$

Замечание

В условиях регулярности, как мы помним, вклад имеет нулевое среднее, поэтому информационная матрица будет его матрицей ковариации:

$$I_{\mathbf{X}}(\theta) = E_{\theta}(U_{\theta}(\mathbf{X})U_{\theta}^T(\mathbf{X})) = D_{\theta}U_{\theta}(\mathbf{X}).$$

- 4) если $\tau : \Theta \rightarrow \mathbb{R}^d$ - гладкая вектор-функция, то через $\tau'(\theta)$ мы обозначим ее матрицу Якоби:

$$\tau'(\theta) = \left(\frac{\partial \tau_i(\theta)}{\partial \theta_j}, i = 1, \dots, d, j = 1, \dots, k \right) \in \mathbf{Mat}(d \times k).$$

- 5) если $\hat{\theta}(\mathbf{X}) \in \mathbb{R}^d$ - то несмещенная оценка $\tau(\theta)$, то ее матрица ковариаций есть

$$D_{\theta}\hat{\theta}(\mathbf{X}) = E_{\theta} \left((\hat{\theta}(\mathbf{X}) - \tau(\theta))(\hat{\theta}(\mathbf{X}) - \tau(\theta))^T \right) \in \mathbf{Mat}(d \times d).$$

Матричное неравенство Коши-Буняковского

Для доказательства многомерного аналога неравенства Рао-Крамера нам понадобится и многомерный аналог неравенства Коши-Буняковского.

Лемма (матричное неравенство Коши-Буняковского) Пусть Ψ и H - случайные матрицы, $\Psi \in \text{Mat}(n' \times m')$, $H \in \text{Mat}(k' \times m')$, причем матрица $E(HH^T)$ обратима. Тогда

$$E(\Psi\Psi^T) \geq E(\Psi H^T)(EHH^T)^{-1}E(H\Psi^T),$$

причем равенство достигается тогда и только тогда, когда $\Psi = ZH$, где

$$Z = E(\Psi H^T)(EHH^T)^{-1}.$$

Замечание

Запись $A \geq B$ означает, что матрица $A - B$ неотрицательно определена.

Доказательство.

Для начала вспомним, что для любой матрицы A матрица $AA^T \geq 0$. Тогда для любой неслучайной матрицы $Z \in \text{Mat}(n' \times k')$ выполнено следующее:

$$(\Psi - ZH)(\Psi - ZH)^T \geq 0.$$

Возьмем математическое ожидание:

$$E(\Psi - ZH)(\Psi - ZH)^T \geq 0.$$

Теперь раскроем по линейности:

$$E\Psi\Psi^T - ZE\Psi H^T - (E\Psi H^T)Z^T + Z(EHH^T)Z^T \geq 0.$$

Положим $Z = E\Psi H^T(EHH^T)^{-1}$. Тогда вышеприведенное неравенство примет нужный вид:

$$E\Psi\Psi^T - E\Psi H^T(EHH^T)^{-1}EHH^T E\Psi^T \geq 0.$$

Для критерия равенства остается заметить, что

$$E(\Psi - ZH)(\Psi - ZH)^T = 0 \Leftrightarrow \Psi = ZH.$$

■

Условия регулярности

Для выполнения неравенства Рао-Крамера снова будут нужны условия регулярности. Они во многом аналогичны одномерному случаю.

- 1) Множество значений параметра $\Theta \subset \mathbb{R}$ - открытое связное множество в \mathbb{R}^k .
- 2) Множество $A = \{x : p_{\theta} > 0\}$ не зависит от θ
- 3) Для любой статистики $S(\mathbf{X})$ с равномерно ограниченным по θ вторым моментом $E_{\theta} S^2(\mathbf{X}) < M < +\infty$ выполнено: для любого $\theta \in \Theta$

$$\frac{\partial}{\partial \theta_j} E_{\theta} S(\mathbf{X}) = E_{\theta} \left(S(\mathbf{X}) \cdot \frac{\partial}{\partial \theta_j} \ln p_{\theta}(\mathbf{X}) \right),$$

$$j = 1, \dots, k.$$

- 4) Информационная матрица $I_{\mathbf{X}}(\theta)$ конечна и положительно определена для всех $\theta \in \Theta$.

Многомерное неравенство Рао-Крамера

Теперь мы готовы сформулировать и доказать многомерное неравенство Рао-Крамера.

Теорема 6.1. (многомерное неравенство Рао-Крамера)

Пусть выполнены условия регулярности. Пусть $\hat{\theta}(\mathbf{X})$ - это несмещенная оценка $\tau(\theta) \in \mathbb{R}^d$ с равномерно ограниченной по θ матрицей $E_{\theta} \hat{\theta}(\mathbf{X}) \hat{\theta}(\mathbf{X})^T$. Тогда для любого $\theta \in \Theta$ выполняется неравенство

$$D_{\theta} \hat{\theta}(\mathbf{X}) \geq \tau'(\theta) I_{\mathbf{X}}^{-1}(\theta) (\tau'(\theta))^T,$$

где $\tau'(\theta)$ - матрица Якоби размера $(d \times k)$

Доказательство.

Доказательство повторяет рассуждения в одномерном случае. Подставляя $S(\mathbf{X}) = 1$ в условие 3, получаем, что

$$0 = E_{\theta} \left(\frac{\partial}{\partial \theta_j} \ln p_{\theta}(\mathbf{X}) \right),$$

$j = 1, \dots, k$, то есть $E_{\theta} U_{\theta}(\mathbf{X}) = 0$. Подставляя $S(\mathbf{X}) = \hat{\theta}_i(\mathbf{X})$, получаем

$$\frac{\partial \tau_i(\theta)}{\partial \theta_j} = E_{\theta} \left(\hat{\theta}_i(\mathbf{X}) \cdot \frac{\partial}{\partial \theta_j} \ln p_{\theta}(\mathbf{X}) \right),$$

или в матричном виде

$$\tau'(\theta) = E_{\theta} \hat{\theta}(\mathbf{X}) (U_{\theta}(\mathbf{X}))^T \stackrel{\text{т.к.}}{=}_{E_{\theta} U_{\theta}(\mathbf{X})=0} E_{\theta} (\hat{\theta}(\mathbf{X}) - \tau(\theta)) (U_{\theta}(\mathbf{X}))^T.$$

Обозначим

$$\Psi = \hat{\theta}(\mathbf{X}) - \tau(\theta) \in \text{Mat}(d \times 1),$$

$$\mathbf{H} = U_{\theta}(\mathbf{X}) \in \text{Mat}(k \times 1).$$

Тогда

$$\tau'(\theta) = E_{\theta} \Psi \mathbf{H}^T.$$

Согласно матричному неравенству Коши-Буняковского, получаем искомое неравенство:

$$D_{\theta} \hat{\theta}(\mathbf{X}) = E \Psi \Psi^T \geq E \Psi \mathbf{H}^T (E \mathbf{H} \mathbf{H}^T)^{-1} E \mathbf{H} \Psi^T = \tau'(\Psi) I_{\mathbf{X}}^{-1}(\Psi) (\tau'(\Psi))^T.$$

■

Критерий равенства

Как и в одномерном случае, можно предложить явный критерий достижения равенства в многомерном неравенстве Рао-Крамера.

Утверждение (критерий равенства)

В многомерном неравенстве Рао-Крамера равенство для всех θ достигается тогда и только тогда, когда для всех $\theta \in \Theta$ с вероятностью 1 выполняется равенство

$$\hat{\theta}(\mathbf{X}) - \tau(\theta) = \tau'(\theta) I_{\mathbf{X}}^{-1}(\theta) U_{\theta}(\mathbf{X}).$$

Доказательство этого утверждения оставляется читателю в качестве упражнения.

О выполнении условий регулярности

Как и в одномерном случае можно предложить достаточное условие выполнения сложного условия регулярности 3.

Теорема 6.2. Пусть для почти всех x по мере μ функция $\sqrt{p_{\theta}(x)}$ непрерывно дифференцируема по θ_j , $j = 1, \dots, k$, а также пусть информационная матрица $I_{\mathbf{X}}(\theta)$ положительно определена и непрерывна по θ . Тогда выполнено условие регулярности 3.

Функция правдоподобия

Предположим, что мы получили некоторые данные, например, вектор из нулей и единиц: 0,0,1,1,1,0,1,1,0,1. Получилось 10 чисел: 6 единиц и 4 нуля.

Попытаемся понять, какое это распределение? У нас есть 2 гипотезы: это биномиальное распределение $\text{Bin}(1, \frac{1}{2})$ (схема Бернулли симметричная), или это пуассоновское распределение $\text{Pois}(\frac{1}{2})$ с параметром $\frac{1}{2}$. Возникает вопрос: какое распределение взять?

В качестве критерия проверки можно предложить следующее: попробуем понять, с какой вероятностью данная нам последовательность реализовалась в каждом распределении.

- Для биномиального - $(\frac{1}{2})^{10}$.
- Для пуассоновского распределения ноль выпадает с вероятностью $e^{-\frac{1}{2}}$, единица с вероятностью $\frac{1}{2}e^{-\frac{1}{2}}$, поэтому мы получаем $(\frac{1}{2})^6 e^{-5}$.

Сравнив полученные 2 числа можно сделать вывод, что для схемы Бернулли выпадение данного набора чисел более вероятно, чем для пуассоновского распределения с параметром $\frac{1}{2}$. Подобный подход сравнения - метод максимального правдоподобия, то есть мы выбираем те распределения, для которых реализация конкретного набора данных более вероятна, чем для других.

Введем понятие функции правдоподобия.

Определение 6.1. Пусть \mathbf{X} - наблюдение с неизвестным распределением $P \in \{P_\theta, \theta \in \Theta\}$, где $\{P_\theta, \theta \in \Theta\}$ есть доминируемое семейство с плотностью $p_\theta(\mathbf{X})$ по мере μ .

Тогда функцией правдоподобия называется случайная величина $f_\theta(\mathbf{X}) = p_\theta(\mathbf{X})$.

Замечание

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

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

Оценка максимального правдоподобия

Определение 6.2. Оценкой параметра θ по методу максимального правдоподобия, или же оценкой максимального правдоподобия (ОМП), называется

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

Данное определение уже накладывает несколько ограничений: во-первых, что максимум существует и единственен, а во-вторых, что он является борелевской функцией от \mathbf{X} .

Философия метода: "мы живем в наиболее вероятном мире". Выпавший результат наблюдения должен был произойти с как можно большей вероятностью.

Этот метод хорошо иллюстрируется, когда мы говорим о дискретных распределениях, так как мы сравниваем вероятности получения конкретного набора данных.

В абсолютно непрерывном случае мы в качестве вероятности получить значение x используем плотность в точке x . Плотность не совсем подходит, так как может быть больше 1, но идея такова.

Пример 1

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

Решение

Когда мы говорим об оценке максимального правдоподобия, важно говорить о том, каково множество значений, то есть по какому множеству мы будем максимизировать нашу функцию плотности. В данном случае мы возьмем все множество. Выпишем функцию правдоподобия:

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

Теперь нам нужно максимизировать ее, как функцию от θ . Заметим, что θ^{-n} - это монотонно убывающая функция, поэтому нужно взять такое минимальное θ , что функция правдоподобия не обратится в нуль. Стало быть, $\hat{\theta}(\mathbf{X}) = X_{(n)}$.

Пример 2

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из нормального распределения $\mathcal{N}(a, \sigma^2)$, $a \in \mathbb{R}$, $\sigma > 0$. Найдите оценку максимального правдоподобия параметра $\boldsymbol{\theta} = (a, \sigma^2)$.

Решение

Опять же, начнем с того, что выпишем функцию правдоподобия:

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

$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - a)^2 \right\}.$$

Нам нужно максимизировать ее одновременно по a и по σ^2 . Но работать с экспонентой достаточно неудобно, поэтому прологарифмируем ее:

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

Это гладкая функция, поэтому приравняем частные производные к нулю и решаем систему (для удобства возьмем параметром σ^2):

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

Отсюда несложно получить, что решением будут оценки

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

Заметим, что решение получилось единственным, это говорит о том, что это будет точка максимума. Это несложно проверить стандартными методами анализа. Тем самым, $\hat{\theta}(\mathbf{X}) = (\bar{X}, S^2)$.

В следующий раз

На следующей лекции мы будем обсуждать асимптотические свойства оценки максимального правдоподобия. Мы увидим, что эти свойства достаточно хороши, и оценка максимального правдоподобия поставляет нам наилучшие оценки в асимптотическом подходе к сравнению оценок.

Оценка максимального правдоподобия хороша с практической точки зрения, так как ее нахождение требует решения экстремальных задач, задач оптимизации, а методы решения этих задач весьма разнообразны и продвинуты.

7. Лекция 7

На этой лекции мы разберем свойства оценки максимального правдоподобия. В прошлый раз мы ввели понятие оценки максимального правдоподобия и разобрали несколько примеров. Сегодня мы разберем, почему эта оценка весьма хороша.

Экстремальное свойство правдоподобия

У оценок максимального правдоподобия и функции правдоподобия есть хорошие асимптотические свойства (состоятельность и асимптотическая нормальность). Но они требуют некоторых условий регулярности. Будем постепенно формулировать их и доказывать свойства.

R0 Параметрическое семейство распределений $\{P_\theta, \theta \in \Theta\}$ - это доминируемое семейство с плотностью $p_\theta(x)$ по мере μ и различными распределениями, то есть $P_{\theta_0} = P_{\theta_1}$ тогда и только тогда, когда $\theta_0 = \theta_1$.

R1 $\mathbf{X} = (X_1, \dots, X_n)$ - выборка растущего размера из неизвестного распределения $P \in \{P_\theta, \theta \in \Theta\}$.

R2 $A = \{x : p_\theta > 0\}$ не зависит от θ . Здесь $p_\theta(x)$ - плотность одного элемента выборки.

Напомним, что через $f_\theta(X_1, \dots, X_n)$ мы обозначаем функцию правдоподобия.

Теорема 7.1. (экстремальное свойство правдоподобия)

В условиях регулярности **R0-R2** для всех различных $\theta_1, \theta_2 \in \Theta$ выполнено

$$\lim_{n \rightarrow \infty} P_{\theta_0}(f_{\theta_0}(X_1, \dots, X_n) > f_{\theta_1}(X_1, \dots, X_n)) = 1.$$

Смысл утверждения: при сравнении значений функций правдоподобия в истинном значении параметра и не истинном, значение в точке истинного значения будет почти всегда больше.

Доказательство.

Будем считать, что все X_i принадлежат A , то есть $p_\theta(x) > 0$. Посмотрим, при каких условиях выполняется событие $f_{\theta_0}(\mathbf{X}) > f_{\theta_1}(\mathbf{X})$. Для этого прологарифмируем и преобразуем выражение:

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

В левой части последнего неравенства стоит сумма независимых одинаково распределенных случайных величин, которая делится на их количество. По усиленному закону больших чисел

$$\frac{1}{n} \sum_{i=1}^n \frac{p_{\theta_1}(X_i)}{p_{\theta_0}(X_i)} < 0 \xrightarrow{P_{\theta_0}-\text{п.н.}} E_{\theta_0} \left(\ln \frac{p_{\theta_1}(X_1)}{p_{\theta_0}(X_1)} \right).$$

Теперь докажем, что

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

Примечание

Если добавить минус к этому математическому ожиданию, то мы 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 всюду.

Логарифм - функция, которая выпукла вверх, поэтому воспользуемся неравенством Йенсена:

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

Мы поняли, что математическое ожидание ≤ 0 . Почему же неравенство, которое необходимо проверить, строгое? Предположим, что это не так, то есть

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

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 всюду. Но $\ln(x)$ нелинейна. Тогда получаем, что аргумент должен быть равен единице почти везде: $\mu(\{x : p_{\theta_0}(x) \neq p_{\theta_1}(x)\}) = 0$. Но это означает, что $P_{\theta_0} = P_{\theta_1}$, что противоречит условию **R0**. В итоге получаем, что

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

Следствие

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 $\theta_0 \in \Theta$

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

Состоятельность решения уравнения правдоподобия

Теперь попробуем разобраться со случаем, когда множество значений параметра не является дискретным множеством. Введем еще два условия регулярности:

R3 Θ есть открытый интервал на \mathbb{R} .

R4 $p_\theta(x)$ непрерывно дифференцируема по θ для всех $a \in A$.

Теорема 7.2. (состоятельность решения уравнения правдоподобия)

В условиях регулярности **R0-R4** уравнение правдоподобия

$$\frac{\partial}{\partial \theta} f_\theta(X_1, \dots, X_n) = 0$$

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

Доказательство.

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

$$A_n = \{f_{\theta_0}(X_1, \dots, X_n) > f_{\theta_0 + \delta}(X_1, \dots, X_n), f_{\theta_0}(X_1, \dots, X_n) > f_{\theta_0 - \delta}(X_1, \dots, X_n)\}$$

Тогда согласно экстремальному свойству правдоподобия

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

Попробуем изобразить событие A_n и понять, как мы можем найти корни уравнения правдоподобия.

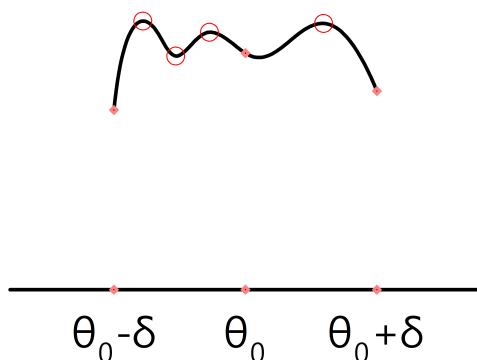


Рис. 4

Рассмотрим наше множество Θ и указанный отрезок. Мы знаем, что значение функции в середине отрезка больше, чем на краях. Функция гладкая, поэтому на отрезке

возникают локальные максимумы, то есть точки, в которых производная обнуляется, соответственно, внутри отрезка мы обязательно находим одно из решений уравнения правдоподобия. Если это событие выполнено, то мы получаем, что решение существует, причем это решение будет в некоторой небольшой окрестности точки θ_0 .

Допустим, что на этом отрезке есть несколько корней (не обязательно конечное число). Пусть $\tilde{\theta}(\mathbf{X})$ - ближайший к θ_0 корень уравнения правдоподобия. Он корректно определен, так как в силу непрерывности производной множество корней замкнуто. Тогда для любого $\varepsilon > 0$ будет выполнено

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

Действительно, зафиксируем ε и заметим, что рассуждения выше верны и 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 тоже лежит на этом отрезке с вероятностью, стремящейся к 1. ■

Мы доказали теорему. Но она, на самом деле, задает больше вопросов, чем дает ответов.

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

Впрочем, если уравнение правдоподобия всегда имеет только один корень, то все становится на свои места. Первый, второй (корень один, поэтому он не зависит от параметра) и четвертый вопросы получают ответы. И остается только третий.

Следствие

Если в условиях **R0-R4** 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, значит, оценка максимального правдоподобия будет состоятельной.

Доказательство.

По сути, доказательство повторяет предыдущие рассуждения. Пусть θ_0 - истинное

значение параметра. Снова возьмем любое $\delta > 0$ такое, что $[\theta_0 - \delta, \theta_0 + \delta] \subset \Theta$ и заметим, что

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

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 $\hat{\theta}(\mathbf{X})$.

Осталось понять, почему это точка глобального максимума. Мы также предполагаем, что значение в точке θ_0 больше, чем на концах отрезка. Точка максимума $\hat{\theta}(\mathbf{X})$ - решение уравнения правдоподобия. Может ли быть так, что $\hat{\theta}(\mathbf{X})$ - не глобальный максимум? Допустим, что есть точка, например, вне нашего отрезка, в котором значение больше. В силу гладкости функции мы получаем, что между этими двумя точками есть точка локального минимума, где производная обнуляется. Поэтому, если наша точка не глобальный максимум, то обязательно найдется еще одно решение уравнения, что противоречит предположению. 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 $\varepsilon > 0$

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

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

Асимптотическая нормальность решения уравнения правдоподобия

В некоторых дополнительных условиях регулярности оценка максимального правдоподобия будет не только состоятельной, но еще и асимптотически нормальной.

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

R6 Интеграл

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

можно дважды дифференцировать под знаком интеграла.

R7 Для всех $\theta \in \Theta$ информация Фишера положительна и конечна

$$0 < i(\theta) = E_{\theta} \left(\frac{\partial}{\partial \theta} \ln p_{\theta}(X_1) \right)^2 < +\infty,$$

где $i(\theta)$ - информация Фишера одного элемента выборки.

R8 Для любого $\theta_0 \in \Theta$ существует $\delta > 0$ и функция $M(x)$ такая, что для всех $\theta \in [\theta_0 - \delta, \theta_0 + \delta]$

$$\left| \frac{\partial^3}{\partial \theta^3} \ln p_{\theta}(x) \right| \leq M(x), \text{ причем } E_{\theta_0} M(X_1) < +\infty.$$

Сформулируем теорему об асимптотической нормальности.

Теорема 7.3. В условиях регулярности **R0-R8** любая состоятельная последовательность $(\hat{\theta}(\mathbf{X}), n \in \mathbb{N})$ корней уравнения правдоподобия удовлетворяет свойству асимптотической нормальности: для всех $\theta_0 \in \Theta$

$$\sqrt{n}(\hat{\theta}(X_1, \dots, X_n) - \theta_0) \xrightarrow{d_{\theta_0}} \mathcal{N}\left(0, \frac{1}{i(\theta_0)}\right).$$

Доказательство.

Обозначим через $\mathcal{L}(\mathbf{X}, \theta) = \ln f_{\theta}(\mathbf{X})$ логарифмическую функцию правдоподобия, и введем обозначение $\mathcal{L}^{(n)}(\mathbf{X}, \theta)$ для n -й частной производной $\mathcal{L}(\mathbf{X}, \theta)$ по θ .

Далее, пусть θ_0 есть истинное значение параметра, то есть $\hat{\theta}(\mathbf{X})$ сходится к θ_0 по вероятности 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 . Далее 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 $\tilde{\theta}_n = \tilde{\theta}_n(\mathbf{X}, \theta_0)$ находится между $\hat{\theta}_n(\mathbf{X})$ и θ_0 . Теперь вспомним, что $\hat{\theta}_n(\mathbf{X})$ есть решения уравнения правдоподобия. Тогда левая часть равенства $\mathcal{L}'(\mathbf{X}, \hat{\theta}_n(\mathbf{X})) = 0$. Теперь выразим $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta_0)$ из данного равенства:

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

Теперь рассмотрим это выражение по частям.

- Начнем с числителя. Заметим, что

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

Как мы уже знаем, возможность дифференцирования под знаком интеграла показывает, что $E_{\theta_0} U_{\theta_0}(X_1) = 0$. Тогда по центральной предельной теореме

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

- Теперь рассмотрим первое слагаемое в знаменателе. По усиленному закону больших чисел

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

Докажем, что

$$-E_{\theta_0} \left(\frac{\partial^2}{\partial \theta^2} \ln p_{\theta}(X_1) \right) = i(\theta).$$

Заметим, что

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

Теперь возьмем математическое ожидание и воспользуемся тем, что мы можем два раза дифференцировать под знаком интеграла:

$$\begin{aligned} E_{\theta_0} \left(\frac{\partial^2}{\partial \theta^2} \ln p_{\theta}(X_1) \right) &= \int_A \frac{\partial^2}{\partial \theta^2} \ln p_{\theta}(x) p_{\theta}(x) \mu(dx) = \\ &= - \int_A \left(\frac{\partial}{\partial \theta} \ln p_{\theta}(x) \right)^2 p_{\theta}(x) \mu(dx) + \int_A \frac{\partial^2}{\partial \theta^2} p_{\theta}(x) \mu(dx) = \\ &= -i(\theta) + \frac{\partial^2}{\partial \theta^2} \int_A p_{\theta}(x) \mu(dx) = -i(\theta). \end{aligned}$$

Отсюда получаем, что

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

- Теперь перейдем ко второму слагаемому в знаменателе и покажем, что он стремится к нулю по вероятности. Заметим, что по условию $\hat{\theta}_n(\mathbf{X})$ сходится к θ_0 по вероятности. Из этого можно сделать вывод, что $\tilde{\theta}_n$ тоже стремится к θ_0 по вероятности. Далее, по **R8**

$$\left| \frac{1}{n} \mathcal{L}'''(\mathbf{X}, \tilde{\theta}_n) \right| \leq \frac{1}{n} \sum_{i=1}^n M(X_i) \xrightarrow{P_{\theta_0}} E_{\theta_0} M(X_1).$$

Тогда получаем произведение ограниченной по вероятности случайной величины на сходящуюся к нулю по вероятности. Произведение сходится к нулю по вероятности.

В итоге

$$\frac{1}{2n} \mathcal{L}'''(\mathbf{X}, \tilde{\theta}_n) (\hat{\theta}_n(\mathbf{X}) - \theta_0) \xrightarrow{P_{\theta_0}} 0.$$

Собирая все вместе и применяя лемму Slutsky, получаем, что

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

■

Сама теорема напрямую не говорит про оценку максимального правдоподобия, она снова апеллирует к некоторому решению уравнения правдоподобия. Сформулируем следствие, которое даст нам асимптотическую нормальность оценки максимального правдоподобия.

Следствие

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

Теорема Бахадура

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

Если в условиях регулярности, похожих на **R0-R8** (в частности, условия для третьей производной $\ln p_\theta(x)$ заменяются на условия для второй), то имеет место следующий асимптотический аналог неравенства Рао-Крамера.

Теорема 7.4. (Бахадур)

Если в некоторых условиях регулярности оценка $\hat{\theta}_n(\mathbf{X})$ является асимптотически нормальной оценкой параметра θ с асимптотической дисперсией $\sigma^2(\theta)$, то $\sigma^2(\theta) \geq i^{-1}(\theta)$ почти всюду по мере Лебега.

Теорему Бахадура мы докажем на следующей лекции, а сейчас обсудим ее результат.

Рассмотрим пример, показывающий, что неравенство действительно может не выполняться в отдельных точках.

Пример

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

$$\hat{\theta}_n(\mathbf{X}) = \begin{cases} \bar{\mathbf{X}}, & \text{если } |\bar{\mathbf{X}}| \geq n^{-\frac{1}{4}}; \\ \beta \bar{\mathbf{X}}, & \text{если } |\bar{\mathbf{X}}| < n^{-\frac{1}{4}}, \end{cases}$$

где $\beta \in (0, 1)$ фиксированное число. Найдите асимптотическую дисперсию $\sigma^2(\theta)$ и сравните ее с обратной информацией Фишера $i^{-1}(\theta)$ одного элемента.

Решение

Как известно, усиленный закон больших чисел имеет скорость сходимости порядка $O(n^{-\frac{1}{2}})$. Следовательно, если $\theta \neq 0$, то

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

Теперь рассмотрим распределение $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta)$. Для любого $x \in \mathbb{R}$ имеем:

$$\begin{aligned} P_{\theta}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x) &= \\ &= P_{\theta}\left(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x, \hat{\theta}_n(\mathbf{X}) = \bar{\mathbf{X}}\right) + P_{\theta}\left(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x, \hat{\theta}_n(\mathbf{X}) = \beta \bar{\mathbf{X}}\right) = \\ &= P_{\theta}(\sqrt{n}(\bar{\mathbf{X}} - \theta) \leq x) + o(1). \end{aligned}$$

Мы знаем, что $P_{\theta}(\sqrt{n}(\bar{\mathbf{X}} - \theta) \leq x) = \Phi(x)$ - функция распределения $\mathcal{N}(0, 1)$, откуда получаем, что

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

Пусть теперь $\theta = 0$. В таком случае, наоборот,

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

Тогда

$$\begin{aligned} P_{\theta}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x) &= \\ &= P_{\theta}\left(\sqrt{n}\bar{\mathbf{X}} \leq x, \hat{\theta}_n(\mathbf{X}) = \bar{\mathbf{X}}\right) + P_{\theta}\left(\sqrt{n}\bar{\mathbf{X}} \leq \frac{x}{\beta}, \hat{\theta}_n(\mathbf{X}) = \beta\bar{\mathbf{X}}\right) = \\ &= P_{\theta}\left(\sqrt{n}\bar{\mathbf{X}} \leq \frac{x}{\beta}\right) + o(1). \end{aligned}$$

Следовательно,

$$\lim_{n \rightarrow \infty} P_{\theta}(\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \leq x) = \Phi\left(\frac{x}{\beta}\right) \Rightarrow \sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \xrightarrow{d_{\theta}} \mathcal{N}(0, \beta^2).$$

В итоге, $\sqrt{n}(\hat{\theta}_n(\mathbf{X}) - \theta) \xrightarrow{d_{\theta}} \mathcal{N}(0, \sigma^2(\theta))$, где $\sigma^2(\theta) = 1$ для всех $\theta \neq 0$ и $\sigma^2(0) = \beta^2 < 1$. Однако информация Фишера всегда равна 1:

$$i(\theta) = E_{\theta}(X_1 - \theta)^2 = D_{\theta}X_1 = 1.$$

Эффективность оценки максимального правдоподобия

Следствие

В условиях регулярности **R0-R8** и условиях предыдущего следствия оценка максимального правдоподобия является наилучшей оценкой θ в асимптотическом подходе в классе асимптотически нормальных оценок с непрерывной асимптотической дисперсией.

Определение 7.1. 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 оценка $\hat{\theta}_n(\mathbf{X})$ называется асимптотически эффективной оценкой θ .

Оказывается, что оценка максимального правдоподобия не только асимптотически эффективна, но и просто эффективна.

Теорема 7.5. (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 при $\tau(\theta) = \theta$: для любого $\theta \in \Theta$ выполнено

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

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

8.

Лекция 8

Теорема Бахадура

В прошлый раз мы узнали, что в условиях регулярности оценка максимального правдоподобия является асимптотически нормальной с асимптотической дисперсией $\frac{1}{i(\theta)}$.

Вопрос: может ли асимптотическая дисперсия быть меньше?

Оказывается, принципиально это невозможно. Данный результат носит название теоремы Бахадура. Для его точной формулировки нам снова понадобятся условия регулярности.

R1 Параметрическое семейство распределений $\{P_\theta, \theta \in \Theta\}$ - это доминируемое семейство с плотностью $p_\theta(x)$ по мере μ и различными распределениями, то есть $P_{\theta_0} = P_{\theta_1}$ тогда и только тогда, когда $\theta_0 = \theta_1$.

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

R3 Θ есть открытый интервал на \mathbb{R} .

R4 Функция $l(x, \theta) = \ln p_\theta(x)$ дважды непрерывно дифференцируема по θ для всех $x \in \mathcal{X}$.

R5 Информация Фишера

$$i(\theta) = E_\theta(l'(X_1, \theta))^2$$

положительна и конечна для всех $\theta \in \Theta$.

R6 Для всех $\theta \in \Theta$ выполнены равенства

$$E_\theta l'(X_1, \theta) = 0, \quad -E_\theta l''(X_1, \theta) = i(\theta),$$

R7 Для любого $\theta_0 \in \Theta$ существуют такие $\delta = \delta(\theta_0) > 0$ и функция $M(x) \geq 0$, что для всех $\theta \in [\theta_0 - \delta, \theta_0 + \delta]$ выполнено

$$|l''(x, \theta)| \leq M(x) \text{ и } E_{\theta_0} M(X_1) < +\infty.$$

Теорема 8.1. (Бахадур)

Если в условиях регулярности **R1-R7** оценка $\tau_n(\mathbf{X})$ является асимптотически нормальной оценкой параметра θ с асимптотической дисперсией $\sigma^2(\theta)$, то $\sigma^2(\theta) \geq \frac{1}{i(\theta)}$ почти всюду по мере Лебега.

Вспомогательная теорема

Сразу приступить к ее доказательству мы не можем, нам понадобится одна вспомогательная теорема, которая представляет самостоятельный интерес.

Теорема 8.2. Пусть в условиях регулярности **R1-R7** для последовательности оценок $\{\tau_n(\mathbf{X}), n \in \mathbb{N}\}$ параметра θ выполнено свойство асимптотической нормальности: для всех $\theta \in \Theta$

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

Если для некоторого $\theta_0 \in \Theta$ и $\theta_n = \theta_0 + n^{-\frac{1}{2}}$ выполнено неравенство

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

$$\text{то } \sigma^2(\theta_0) \geq \frac{1}{i(\theta_0)}.$$

Доказательство.

Для доказательства Теоремы 8.2 необходимо сперва доказать несколько вспомогательных утверждений и ввести несколько обозначений. Начнем с того, что введем логарифмическую функцию правдоподобия:

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

Далее, введем следующую случайную величину:

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

Покажем, что для случайной величины T_n выполнено одно интересное свойство.

Лемма (1)

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 $l(x, \theta)$ и разложим $L(\mathbf{X}, \theta_n)$ в ряд Тейлора в окрестности θ_0 :

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

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

где $\tilde{\theta}_n \in (\theta_0, \theta_n)$ - некоторая промежуточная точка. Далее, мы хотим понять, куда сходятся слагаемые, так как если перенести из правой части равенства $\mathcal{L}(\mathbf{X}, \theta_0)$ в левую, то мы получим выражение $\mathcal{L}(\mathbf{X}, \theta_n) - \mathcal{L}(\mathbf{X}, \theta_0)$, которое стоит в скобках в выражении для T_n . Введем следующую случайную величину:

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

Докажем, что эта последовательность стремится к нулю почти наверное:

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

Для этого введем следующую функцию: для $\varepsilon \in (0, \delta(\theta_0))$ ($\delta(\theta_0)$ - величина из условия **R7**)

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

Далее, положим $m(\varepsilon) = E_{\theta_0} A(X_1, \varepsilon)$. Заметим, что $A(x, \varepsilon) \downarrow 0$ при $\varepsilon \rightarrow 0+$, так как чем меньше ε , тем меньше множество, по которому мы берем \sup . Далее, воспользуемся свойством регулярности **R7**:

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

Тем самым, можно применить теорему Лебега о мажорируемой сходимости. Из нее следует, что $m(\varepsilon) \rightarrow 0$ при $\varepsilon \rightarrow 0$.

Далее, заметим, что если взять $\varepsilon^2 > \frac{1}{n}$, то

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

Однако по усиленному закону больших чисел

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

Из вышеуказанного следует, что для любого $\varepsilon > 0$

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

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

Теперь преобразуем T_n следующим образом:

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

Осталось найти предел по распределению данного выражения.

- Начнем с $n^{-\frac{1}{2}} \mathcal{L}'(\mathbf{X}, \theta_0)$. Заметим, что из условия регулярности **R6** и центральной предельной теоремы следует, что

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

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

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

- Из доказанного выше следует, что

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

В итоге, по лемме Slutsky получаем, что

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

Лемма (1) доказана. ■

Простым следствием леммы (1) является следующее утверждение.

Лемма (2)

Для любых $y, \alpha > 0$ выполнено

$$E_{\theta_0}(e^{\alpha T_n} \cdot I\{T_n < y\}) \rightarrow E(e^{\alpha \xi} \cdot I\{\xi < y\})$$

при $n \rightarrow +\infty$, где $\xi \sim \mathcal{N}(0, 1)$.

Доказательство.

Рассмотрим функцию $f(x) = e^{\alpha \min(x,y)}$. Она является непрерывной и ограниченной на \mathbb{R} . Тогда в силу леммы 1

$$E_{\theta_0} f(T_n) \rightarrow E f(\xi).$$

С другой стороны,

$$f(x) - e^{\alpha x} \cdot I\{x < y\} = e^{\alpha y} \cdot I\{x \geq y\}.$$

Стало быть,

$$\begin{aligned} E_{\theta_0}(e^{\alpha T_n} \cdot I\{T_n < y\}) &= E_{\theta_0} f(T_n) - E_{\theta_0}(e^{\alpha y} \cdot I\{T_n \geq y\}) = E_{\theta_0} f(T_n) - e^{\alpha y} P_{\theta_0}(T_n \geq y) \rightarrow \\ &\rightarrow E f(\xi) - e^{\alpha y} P(\xi \geq y) = E f(\xi) - E(e^{\alpha y} \cdot I\{\xi \geq y\}) = E(e^{\alpha \xi} \cdot I\{\xi < y\}). \end{aligned}$$

Лемма (2) доказана. ■

Лемма (3)

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

$$\lim_{n \rightarrow \infty} P_{\theta_n}(T_n \geq y) = 1 - \Phi(y - \sqrt{i(\theta_0)}).$$

Доказательство.

Напомним, что T_n есть функция от выборки: $T_n = T_n(\mathbf{X})$. Рассмотрим вероятность

$$\begin{aligned} P_{\theta_n}(T_n < y) &= \int_{\mathcal{X}^n} p_{\theta_n}(x_1, \dots, x_n) I\{T_n(x_1, \dots, x_n) < y\} \mu(dx_1, \dots, dx_n) = \\ &= \int_{\mathcal{X}^n} e^{\mathcal{L}(\mathbf{x}, \theta_n)} I\{T_n(\mathbf{x}) < y\} \mu(d\mathbf{x}). \end{aligned}$$

Вспомним, что из определения статистики T_n следует, что

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

Подставим это выражение в интеграл. Тогда по лемме (2)

$$\begin{aligned} P_{\theta_n}(T_n < y) &= e^{-\frac{i(\theta_0)}{2}} \int_{\mathcal{X}^n} e^{\mathcal{L}(\mathbf{x}, \theta_0) + T_n(\mathbf{x}) \sqrt{i(\theta_0)}} I\{T_n(\mathbf{x}) < y\} \mu(d\mathbf{x}) = \\ &= e^{-\frac{i(\theta_0)}{2}} E_{\theta_0} \left(e^{T_n(\mathbf{X}) \sqrt{i(\theta_0)}} I\{T_n(\mathbf{X}) < y\} \right) \rightarrow e^{-\frac{i(\theta_0)}{2}} E \left(e^{\sqrt{i(\theta_0)} \cdot \xi} I\{\xi < y\} \right), \end{aligned}$$

где $\xi \sim \mathcal{N}(0, 1)$. Осталось заметить, что предел равен искомому выражению:

$$e^{-\frac{i(\theta_0)}{2}} E \left(e^{\sqrt{i(\theta_0)} \cdot \xi} I\{\xi < y\} \right) = e^{-\frac{i(\theta_0)}{2}} \int_{-\infty}^y e^{\sqrt{i(\theta_0)} \cdot z} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz =$$

$$= \int_{-\infty}^y \frac{1}{\sqrt{2\pi}} e^{-\frac{(z - \sqrt{i(\theta_0)})^2}{2}} dz = \Phi(y - \sqrt{i(\theta_0)}).$$

Лемма (3) доказана. ■

Мы близки к завершению доказательства теоремы 8.2. Зафиксируем $y > \sqrt{i(\theta_0)}$ и введем два события:

$$D_n = \{\tau_n \geq \theta_n\} \text{ и } S_n = \{T_n \geq y\}.$$

Согласно лемме 3 получаем, что

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

Теперь заметим, что по условию (1) теоремы 8.2:

$$\frac{1}{2} \geq \lim_{n \rightarrow \infty} P_{\theta_n}(\tau_n(\mathbf{X}) \leq \theta_n) = 1 - \overline{\lim}_{n \rightarrow \infty} P_{\theta_n}(D_n).$$

Следовательно, $\overline{\lim}_{n \rightarrow \infty} P_{\theta_n}(D_n) \geq \frac{1}{2}$. Значит, существует некоторая подпоследовательность индексов $\{n_k, k \in \mathbb{N}\}$ такая, что

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

Это означает, что при достаточно больших k выполнено

$$P_{\theta_{n_k}}(D_{n_k}) > P_{\theta_{n_k}}(S_{n_k}).$$

Следующая лемма говорит, что данное соотношение верно и для θ_0 .

Лемма (4)

Для всех достаточно больших k выполнено

$$P_{\theta_0}(D_{n_k}) > P_{\theta_0}(S_{n_k}).$$

Доказательство.

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

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

где λ - некоторое положительное число, $y' = y\sqrt{i(\theta_0)} - \frac{i(\theta_0)}{2}$.

Тогда для всех \mathbf{x} будет выполнено неравенство:

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

Действительно,

- либо $\mathbf{x} \in S_n$, тогда левая часть равна $p_{\theta_n}(\mathbf{x}) - \lambda p_{\theta_0}(\mathbf{x}) \geq 0$, а индикатор $I\{\mathbf{x} \in D_n\} \leq 1$;
- либо $\mathbf{x} \notin S_n$, тогда левая часть равна нулю, правая неположительна, так как $p_{\theta_n}(\mathbf{x}) - \lambda p_{\theta_0}(\mathbf{x}) < 0$.

Проинтегрируем данное неравенство по мере μ и получим:

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

Простой перегруппировкой мы получим желаемое:

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

Лемма (4) доказана. ■

Теперь найдем пределы $P_{\theta_0}(D_n)$ и $P_{\theta_0}(S_n)$ с ростом n . Согласно лемме (1) мы уже знаем, что

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

Далее, по условию теоремы

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

Следовательно, из леммы (4) мы получаем, что для любого $y > \sqrt{i(\theta_0)}$ выполняется неравенство

$$\frac{1}{\sigma(\theta_0)} \leq y \iff \sigma^2(\theta_0) \geq \frac{1}{y^2}.$$

Отсюда и вытекает искомое неравенство $\sigma^2(\theta_0) \geq \frac{1}{i(\theta_0)}$.

Теорема 8.2 доказана. ■

Доказательство теоремы Бахадура

Приступим к доказательству теоремы Бахадура.

Доказательство.

В силу теоремы 2 нам достаточно проверить, что неравенство (1) выполнено почти всюду по мере Лебега на Θ .

Рассмотрим следующую функцию: для всех $\theta \in \Theta$ положим

$$h_n(\theta) = \left| P_\theta(\tau_n(\mathbf{X}) \leq \theta) - \frac{1}{2} \right| = \left| P_\theta(\sqrt{n}(\tau_n(\mathbf{X}) - \theta) \leq 0) - \frac{1}{2} \right|,$$

и положим $h_n(\theta) = 0$ для $\theta \in \mathbb{R} \setminus \Theta$. Тогда $h_n(\theta) \in [0, \frac{1}{2}]$, кроме того из условия асимптотической нормальности мы получаем, что для всех $\theta \in \mathbb{R}$ выполнено

$$\lim_{n \rightarrow +\infty} h_n(\theta) = 0.$$

Теперь введем вторую функцию $g_n(\theta) = h_n(\theta + n^{-\frac{1}{2}})$, а также рассмотрим случайную величину $\xi \sim \mathcal{N}(0, 1)$. Посчитаем математическое ожидание $g_n(\xi)$.

$$Eg_n(\xi) = \int_{-\infty}^{+\infty} g_n(x) \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_{-\infty}^{+\infty} h_n\left(x + \frac{1}{\sqrt{n}}\right) \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx.$$

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

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

В силу свойств функций h_n выполнено

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

По теореме Лебега о мажорируемой сходимости получаем, что $Eg_n(\xi) \rightarrow 0$. Так как $g_n(x) \in [0, \frac{1}{2}]$, то это означает, что $g_n(\xi)$ сходится к нулю по вероятности. Но тогда можно выбрать подпоследовательность индексов $\{n_k, k \in \mathbb{N}\}$ такую, что $g_{n_k} \rightarrow 0$ п.н.

Остается заметить, что если $g_{n_k} \rightarrow 0$, то для такого θ будет выполнено условие (1):

$$\lim_{k \rightarrow +\infty} P_{\theta_{n_k}}(\tau_{n_k} < \theta_{n_k}) = \frac{1}{2}.$$

В силу доказанного $g_{n_k} \rightarrow 0$ почти всюду по мере Q , где Q - стандартное нормальное распределение. Но это эквивалентно тому, что $g_{n_k}(\theta) \rightarrow 0$ почти всюду по мере Лебега на \mathbb{R} . ■

9. Лекция 9

Эта лекция будет посвящена разделу теории вероятностей, который называется условное математическое ожидание. С помощью условного математического ожидания мы планируем разобраться с теми подходами к сравнению оценок, которые мы еще не разобрали, а именно с равномерным, байесовским и минимаксным.

Условное математическое ожидание

Вопрос: в теории вероятностей есть понятие условной вероятности. Пусть ξ и η - две случайные величины. Как можно определить $E(\xi|\eta = y)$ и $E(\xi|\eta)$?

Рассмотрим самый простой случай - простые случайные величины. Как известно, для них

$$E\xi = \sum_{i=1}^n x_i P(\xi = x_i).$$

Естественно ввести $E(\xi|\eta = y)$ следующим образом:

$$E(\xi|\eta = y) = \sum_{i=1}^n x_i P(\xi = x_i | \eta = y).$$

Вопрос: как обобщить данное определение на случай непростой случайной величины ξ ?

Домножим на $P(\eta = y)$ правую часть равенства. Тогда

$$\sum_{i=1}^n x_i P(\xi = x_i, \eta = y) = E\left(\sum_{i=1}^n x_i \cdot I\{\xi = x_i\} I\{\eta = y\}\right) = E(\xi \cdot I\{\eta = y\}).$$

С другой стороны, если положить $E(\xi|\eta = y) = \phi(y)$, то получим равенство

$$\phi(y)P(\eta = y) = E(\phi(y)I\{\eta = y\}) = E(\phi(\eta)I\{\eta = y\}).$$

Тем самым, разумно сказать, что $E(\xi|\eta = y)$ - это такая функция $\phi(y)$, что для любого y выполнено равенство

$$E(\xi \cdot I\{\eta = y\}) = E(\phi(\eta)I\{\eta = y\}).$$

Вопрос: как обобщить данное определение на случай непростой случайной величины η ?

Пусть (Ω, \mathcal{F}, P) - вероятностное пространство, ξ - случайная величина на нем, а $\mathcal{C} \subset \mathcal{F}$ под- σ -алгебра в \mathcal{F} .

Определение 9.1. Напомним, что σ -алгеброй, порожденной случайной величиной ξ называется

$$\mathcal{F}_\xi = \{\xi^{-1}(B) = \{\xi \in B\} : B \in \mathcal{B}(\mathcal{R})\}.$$

Случайная величина ξ называется \mathfrak{C} -измеримой, если порожденная ею σ -алгебра входит в \mathfrak{C} :

$$\mathcal{F}_\xi \subset \mathfrak{C}.$$

Определение 9.2. Условным математическим ожиданием ξ относительно σ -алгебры \mathfrak{C} называется случайная величина $E(\xi|\mathfrak{C})$, удовлетворяющая следующим свойствам:

- 1) $E(\xi|\mathfrak{C})$ является \mathfrak{C} -измеримой случайной величиной (свойство измеримости);
- 2) для любого $A \in \mathfrak{C}$ выполняется равенство $E(\xi \cdot I_A) = E(E(\xi|\mathfrak{C})I_A)$ или

$$\int_A \xi dP = \int_A E(\xi|\mathfrak{C}) dP \quad (\text{интегральное свойство}).$$

Вопрос: почему $E(\xi|\mathfrak{C}) \neq \xi$ в общем случае?

Ответ: ξ не обязательно является \mathfrak{C} -измеримой! Соответственно, мы пытаемся подобрать такую \mathfrak{C} -измеримую случайную величину, для которой интегрирование по любому событию из \mathfrak{C} дает то же самое, что и для ξ .

Существование УМО

Определение 9.3. Пусть (Ω, \mathcal{F}, P) - вероятностное пространство.

Функция $\nu : \mathcal{F} \rightarrow \mathbb{R}$ называется (конечным) зарядом (или мерой со знаком), если ν - счетно-аддитивна на \mathcal{F} , то есть

$$\nu\left(\bigsqcup_{n=1}^{\infty} A_n\right) = \sum_{i=1}^{\infty} \nu(A_n),$$

где ряд в правой части сходится абсолютно, и

$$\sup_{A \in \mathcal{F}} |\nu(A)| < +\infty.$$

Заряд ν является абсолютно непрерывным относительно P , если из того, что $P(A) = 0$ следует, что $\nu(A) = 0$.

Теорема 9.1. (Радон-Никодим)

Пусть (Ω, \mathcal{F}, P) - вероятностное пространство, а ν - заряд, абсолютно непрерывный относительно P . Тогда существует единственная (с точностью до равенства п.н.)

случайная величина $\eta \in L^1(\Omega, \mathcal{F}, P)$ (с конечным математическим ожиданием) такая, что для любого $A \in \mathcal{F}$ выполнено

$$\nu(A) = \int_A \eta dP = E(\eta \cdot I_A).$$

Замечание

В этом случае $\eta = \frac{d\nu}{dP}$ называется производной Радона-Никодима.

Лемма (о существовании УМО)

Если $E|\xi| < \infty$, то для любой σ -алгебры $\mathfrak{C} \subset \mathcal{F}$ условное математическое ожидание $E(\xi|\mathfrak{C})$ существует и единственно с точностью до равенства почти наверное.

Доказательство.

Рассмотрим вероятностное пространство $(\Omega, \mathfrak{C}, P)$. Для любого $A \in \mathfrak{C}$ положим

$$\nu(A) = \int_A \xi dP = E(\xi \cdot I_A).$$

Это заряд на $(\Omega, \mathfrak{C}, P)$, абсолютно непрерывный относительно P . По теореме Радона-Никодима существует единственная (п.н.) случайная величина $\eta \in L^1(\Omega, \mathfrak{C}, P)$ такая, что для любого $A \in \mathfrak{C}$

$$\nu(A) = \int_A \eta dP = E(\eta \cdot I_A).$$

Случайная величина η является \mathfrak{C} -измеримой и удовлетворяет интегральному свойству. Значит, по определению $\eta = E(\xi|\mathfrak{C})$.

Единственность следует из единственности производной Радона-Никодима в теореме. ■

Замечание

Условное математическое ожидание определяется с точностью до равенства п.н. Иногда бывает важным осуществлять правильный подбор его "варианта".

Мы разобрались со существованием. Оставшуюся часть лекции будем обсуждать свойства УМО.

Дискретные σ -алгебры

Начнем с явной формулы, по которой можно будет вычислить УМО в ситуации, когда \mathfrak{C} - дискретная σ -алгебра.

Напоминание

Дискретная σ -алгебра \mathfrak{C} - это σ -алгебра, которая порождена некоторым (не более чем счетным) разбиением $\{D_n, n \in \mathbb{N}\}$ пространства Ω . То есть, в такую σ -алгебру входят сами D_n и всевозможные их объединения.

Лемма

Если σ -алгебра порождена разбиением $\{D_n, n \in \mathbb{N}\}$ с условием $P(D_n) > 0$ для всех $n \in \mathbb{N}$, то для любой случайной величины ξ с условием $E|\xi| < +\infty$ выполнено

$$E(\xi|\mathfrak{C})(\omega) = \sum_{i=1}^{\infty} \frac{E(\xi \cdot I_{D_n})}{P(D_n)} \cdot I_{D_n}(\omega).$$

Попробуем осознать некоторый смысл УМО. Пусть $\Omega = [0, 1]$ разрезан на 4 части D_i , $i = 1, 2, 3, 4$. Случайная величина ξ на Ω - функция на отрезке. Что же такое условное математическое ожидание? Из последней формулы мы поняли, что в данной ситуации условным математическим ожиданием ξ относительно σ -алгебры, порожденной таким разбиением, будет случайная величина, которая является константой на каждом из множеств D_i .

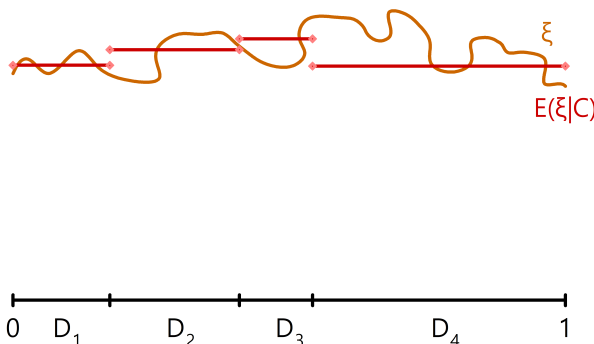


Рис. 5

По сути, условное математическое ожидание - это усреднение нашей случайной величины по множествам из σ -алгебры. В данном случае это происходит явным образом, так как множества простые (не пересекаются, их счетное количество). В общей ситуации все не так явно.

Перейдем теперь к доказательству леммы.

Доказательство.

План доказательства прост. В правой части нам предлагают кандидата в УМО, в силу единственности надо проверить для него два свойства: свойство измеримости и интегральное свойство.

Правая часть есть линейная комбинация несовместных индикаторов событий $D_n \in \mathfrak{C}$, так что это \mathfrak{C} -измеримая случайная величина.

Проверим интегральное свойство. Если $A \in \mathfrak{C}$, то оно представляет собой лишь дизъюнктивное объединение элементов D_n . В связи с этим достаточно рассмотреть ситуацию $A = D_m$ для некоторого $m \in \mathbb{N}$. Имеем:

$$E\left(\left(\sum_{n=1}^{\infty} \frac{E(\xi \cdot I_{D_n})}{P(D_n)} \cdot I_{D_n}\right) I_{D_m}\right) \stackrel{(*)}{=} E\left(\frac{E(\xi \cdot I_{D_m})}{P(D_m)} \cdot I_{D_m}\right) = \frac{E(\xi \cdot I_{D_m})}{P(D_m)} \cdot P(D_m) = E(\xi \cdot I_{D_m}),$$

(*) - из-за несовместности событий D_n остается только слагаемое при $n = m$. ■

Свойства УМО

В общем случае таких хороших формул, как в прошлой лемме, нет, и учиться считать УМО в общем случае мы будем на следующей лекции. Сейчас разберем свойства УМО.

Всюду далее предполагаем конечность необходимых математических ожиданий у рассматриваемых случайных величин, чтобы УМО были корректно определены и существовали.

Свойство (1)

Если случайная величина ξ является \mathfrak{C} -измеримой, то $E(\xi|\mathfrak{C}) = \xi$.

Доказательство.

Очевидно, что для ξ выполняются свойство измеримости и интегральное свойство. ■

Свойство (2, линейность)

Для любых случайных величин ξ , η и σ -алгебры \mathfrak{C} выполнено, что для любых $a, b \in \mathbb{R}$:

$$E(a\xi + b\eta|\mathfrak{C}) = aE(\xi|\mathfrak{C}) + bE(\eta|\mathfrak{C}).$$

Доказательство.

Для начала заметим, что и $E(\xi|\mathfrak{C})$ и $E(\eta|\mathfrak{C})$ являются \mathfrak{C} -измеримыми случайными величинами. Тогда $aE(\xi|\mathfrak{C}) + bE(\eta|\mathfrak{C})$ - это тоже \mathfrak{C} -измеримая случайная величина.

Осталось проверить интегральное свойство. В силу линейности математического ожидания для любого $A \in \mathfrak{C}$

$$E((a\xi + b\eta) \cdot I_A) = aE(\xi \cdot I_A) + bE(\eta \cdot I_A).$$

Согласно интегральному свойству

$$aE(\xi \cdot I_A) + bE(\eta \cdot I_A) = aE(E(\xi|\mathfrak{C})I_A) + bE(E(\eta|\mathfrak{C})I_A)$$

Свернем сумму назад:

$$aE(E(\xi|\mathfrak{C})I_A) + bE(E(\eta|\mathfrak{C})I_A) = E((aE(\xi|\mathfrak{C}) + bE(\eta|\mathfrak{C})) \cdot I_A).$$

Тем самым, получаем желаемое. ■

Свойство (3, формула полной вероятности)

$$E(E(\xi|\mathfrak{C})) = E\xi.$$

Доказательство.

Воспользуемся интегральным свойством, положив $\mathfrak{C} \ni A = \Omega$:

$$E\xi = E(\xi \cdot I_A) = E(E(\xi|\mathfrak{C}) \cdot I_A)E(E\xi|\mathfrak{C})).$$
■

Свойство (4)

Если ξ независима с \mathfrak{C} , то $E(\xi|\mathfrak{C}) = E\xi$.

Доказательство.

Напомним, что ξ независима с \mathfrak{C} тогда и только тогда, когда для любого $A \in \mathfrak{C}$ случайная величина ξ независима с I_A . Заметим, что $E\xi$ - это константа. Тогда $E\xi$ является \mathfrak{C} -измеримой. Проверим интегральное свойство: для любого $A \in \mathfrak{C}$

$$E(\xi \cdot I_A) = |\text{независимость}| = E\xi \cdot E(I_A) = E(E\xi \cdot I_A).$$
■

Свойство (5, сохранение отношения порядка)

Если $\xi \leq \eta$ п.н., то

$$E(\xi|\mathfrak{C}) \leq E(\eta|\mathfrak{C}) \text{ п.н.}$$

Доказательство.

Рассмотрим случайную величину $\delta = E(\xi|\mathfrak{C}) - E(\eta|\mathfrak{C})$. Она является \mathfrak{C} -измеримой

случайной величиной. Теперь рассмотрим $E(\delta \cdot I\{\delta > 0\})$. Заметим, что $\delta \cdot I\{\delta > 0\} \geq 0$. Далее, рассмотрим событие $B = \{\delta > 0\} \in \mathfrak{C}$. Тогда в силу выбора δ

$$E(\delta \cdot I\{\delta > 0\}) = E(E(\xi|\mathfrak{C}) - E(\eta|\mathfrak{C}))I_B).$$

Пользуясь линейностью математического ожидания и интегральным свойством, получаем, что

$$E(\delta \cdot I\{\delta > 0\}) = E(\xi \cdot I_B - \eta \cdot I_B) \leq 0 \text{ так как } \xi \leq \eta.$$

Тогда получаем, что $\delta \cdot I\{\delta < 0\} = 0$ п.н., то есть $\delta \leq 0$ п.н. ■

Свойство (6, неравенство модуля)

$$|E(\xi|\mathfrak{C})| \leq E(|\xi||\mathfrak{C}).$$

Доказательство. Следует из свойств (2) и (5). ■

Свойство (7, телескопическое)

Пусть ξ - случайная величина, а $\mathfrak{C}_1 \subset \mathfrak{C}_2$ - две σ -алгебры. Тогда

$$1) E(E(\xi|\mathfrak{C}_1)|\mathfrak{C}_2) = E(\xi|\mathfrak{C}_1),$$

$$2) E(E(\xi|\mathfrak{C}_2)|\mathfrak{C}_1) = E(\xi|\mathfrak{C}_1).$$

Доказательство.

Первый пункт почти очевиден. Заметим, что случайная величина $E(\xi|\mathfrak{C}_1)$ является \mathfrak{C}_1 измеримой, а, значит, и \mathfrak{C}_2 -измеримой (т.к. $\mathfrak{C}_1 \subset \mathfrak{C}_2$). Далее, применяем свойство (1).

Теперь докажем второй утверждение. Заметим, что свойство измеримости у нас уже есть. Проверим интегральное свойство. Для любого $A \in \mathfrak{C}_1 \subset \mathfrak{C}_2$:

$$\begin{aligned} E(E(\xi|\mathfrak{C}_1)I_A) &= |\text{инт. свойство для } \mathfrak{C}_1| = E(\xi \cdot I_A) = \\ &= |\text{инт. свойство для } \mathfrak{C}_2| = E(E(\xi|\mathfrak{C}_2)I_A). \end{aligned}$$

■

Свойство (8, предельный переход под знаком УМО)

Пусть $\{\xi_n, n \in \mathbb{N}\}$, ξ - случайные величины, $E|\xi| < \infty$, $E|\xi_n| < \infty$. Тогда

$$1) \text{ Если } 0 \leq \xi_n \uparrow \xi \text{ п.н., то } E(\xi_n|\mathfrak{C}) \uparrow E(\xi|\mathfrak{C}) \text{ п.н.}$$

$$2) \text{ Пусть } \xi_n \xrightarrow{\text{п.н.}} \xi \text{ и } |\xi_n| \leq \delta \text{ для всех } n, \text{ причем } E\delta < +\infty. \text{ Тогда}$$

$$E(\xi_n|\mathfrak{C}) \xrightarrow{\text{п.н.}} E(\xi|\mathfrak{C}).$$

Доказательство.

- 1) В силу того, что ξ_n монотонно возрастают, то и $E(\xi_n|\mathfrak{C})$ монотонно возрастают (по свойству (5)). Тогда п.н. существует предел $\eta = \lim_{n \rightarrow \infty} E(\xi_n|\mathfrak{C})$. Проверим, что $\eta = E(\xi|\mathfrak{C})$.

Заметим, что η является \mathfrak{C} -измеримой случайной величиной, как предел п.н. \mathfrak{C} -измеримых случайных величин. Проверим интегральное свойство.

Для любого $A \in \mathfrak{C}$ п.н. выполнено $\xi_n \cdot I_A \uparrow \xi \cdot I_A$ и $E(\xi_n|\mathfrak{C}) \cdot I_A \uparrow \eta \cdot I_A$. Тогда по теореме о монотонной сходимости:

$$E(\xi \cdot I_A) = \lim_{n \rightarrow \infty} E(\xi_n \cdot I_A) = |\text{инт. свойство}| = \lim_{n \rightarrow \infty} E(E(\xi_n|\mathfrak{C})I_A) = E(\eta \cdot I_A).$$

- 2) Рассмотрим $\eta_n = \sup_{m \geq n} |\xi_m - \xi_n|$. Тогда в силу свойства (5)

$$|E(\xi_n|\mathfrak{C}) - E(\xi|\mathfrak{C})| = |E(\xi_n - \xi|\mathfrak{C})| \leq E(|\xi_n - \xi||\mathfrak{C}) \leq E(\eta_n|\mathfrak{C}).$$

Далее, $\eta_n \downarrow 0$ п.н. и $|\eta_n| \leq 2\delta$. Согласно пункту 1) получаем, что

$$E(\eta_n|\mathfrak{C}) \xrightarrow{\text{п.н.}} E(0|\mathfrak{C}) = 0.$$

■

Свойство (9)

Пусть η - это \mathfrak{C} -измеримая случайная величина, а математические ожидания $E\xi\eta$ и $E\xi$ конечны. Тогда

$$E(\xi \cdot \eta|\mathfrak{C}) = \eta E(\xi|\mathfrak{C}).$$

Доказательство.

Заметим, что $\eta E(\xi|\mathfrak{C})$ есть \mathfrak{C} -измеримая случайная величина (произведение двух \mathfrak{C} -измеримых). Проверим интегральное свойство.

Пусть сначала $\eta = I_B$ для некоторого $B \in \mathfrak{C}$. Тогда для любого $A \in \mathfrak{C}$

$$E(\eta\xi \cdot I_A) = E(\xi \cdot I_B \cdot I_A) = E(\xi \cdot I_{A \cap B}) = E(E(\xi|\mathfrak{C})I_{A \cap B}) = E(E(\xi|\mathfrak{C})I_B \cdot I_A) = E(\eta E(\xi|\mathfrak{C})I_A).$$

Заметим, что $E(\eta\xi \cdot I_A)$ и $E(\eta E(\xi|\mathfrak{C})I_A)$ - линейные функции от случайной величины η . Тогда в силу линейности математического ожидания, интегральное свойство будет выполнено и для простых случайных величин вида $\eta = \sum_{k=1}^m c_k I_{B_k}$, $c_k \in \mathbb{R}$, $B_k \in \mathfrak{C}$.

Для произвольной η возьмем последовательность простых \mathfrak{C} -измеримых случайных величин η_n такую, что $\eta_n \xrightarrow{\text{п.н.}} \eta$ и $|\eta_n| \leq |\eta|$. Тогда по свойству (8) получаем, что

$$E(\eta_n \xi|\mathfrak{C}) \xrightarrow{\text{п.н.}} E(\eta \xi|\mathfrak{C}).$$

Однако

$$E(\eta_n \xi | \mathfrak{C}) = \eta_n E(\xi | \mathfrak{C}) \xrightarrow{п.н.} \eta E(\xi | \mathfrak{C}).$$

Следовательно, $E(\xi \eta | \mathfrak{C}) = \eta E(\xi | \mathfrak{C})$. ■

Свойство (10, неравенство Йенсена)

Пусть φ - выпуклая книзу функция. тогда

$$E(\varphi(\xi) | \mathfrak{C}) \geq \varphi(E(\xi | \mathfrak{C})).$$

Доказательство.

В силу выпуклости для любого $x \in \mathbb{R}$ существует $\lambda(x)$ такая, что для любого $y \in \mathbb{R}$ выполнено

$$\varphi(y) \geq \varphi(x) + \lambda(x)(y - x).$$

Положим $y = \xi$, $x = E(\xi | \mathfrak{C})$ и возьмем УМО относительно \mathfrak{C} от обеих частей неравенства. Пользуясь свойствами (1) и (9), получим:

$$\begin{aligned} E(\varphi(\xi) | \mathfrak{C}) &\geq E\left(\varphi(E(\xi | \mathfrak{C})) | \mathfrak{C}\right) + E\left(\lambda(E(\xi | \mathfrak{C}))(\xi - E(\xi | \mathfrak{C})) | \mathfrak{C}\right) = \\ &= |\text{выносим } \lambda(E(\xi | \mathfrak{C})) \text{ как } \mathfrak{C}\text{-измеримую с.в.}| = \\ &= \varphi(E(\xi | \mathfrak{C})) + \lambda(E(\xi | \mathfrak{C}))E(\xi - E(\xi | \mathfrak{C}) | \mathfrak{C}) = \varphi(E(\xi | \mathfrak{C})). \end{aligned}$$

■

10.

Лекция 10

Условное математическое ожидание

Нам понадобятся еще несколько определений, связанных с УМО. Пусть (Ω, \mathcal{F}, P) - вероятностное пространство.

Определение 10.1. Пусть $A \in \mathcal{F}$ - событие. Тогда для под- σ -алгебры $\mathfrak{C} \subset \mathcal{F}$ определим

$$P(A|\mathfrak{C}) := E(I_A|\mathfrak{C}).$$

Если ξ и η - две случайные величины, то полагаем

$$E(\xi|\eta) := E(\xi|\mathcal{F}_\eta),$$

где \mathcal{F}_η - σ -алгебра, порожденная случайной величиной η .

Замечание

Случайная величина ζ является \mathcal{F}_η -измеримой тогда и только тогда, когда существует такая борелевская функция $\phi(x)$, что $\zeta = \phi(\eta)$.

Определение 10.2. Пусть ξ и η - две случайные величины. Тогда величиной $E(\xi|\eta = y)$ называется такая борелевская функция $\varphi(y)$, что для любого $B \in \mathcal{B}(\mathbb{R})$ выполняется равенство

$$E(\xi \cdot I\{\eta \in B\}) = E(\varphi(\eta) \cdot I\{\eta \in B\}) = \int_B \varphi(y) P_\eta(dy).$$

Нам необходима такая конструкция, так как мы хотим вычислять УМО ξ относительно другой случайной величины. Для этого удобно сначала вычислять подобные математические ожидания $E(\xi|\eta = y)$, так как мы имеем дело не со случайными величинами, а с формально более простыми вещами, например, с борелевскими функциями.

Лемма

Если $E|\xi| < +\infty$, то $E(\xi|\eta = y)$ существует и единственная (с точностью до равенства P_η -п.н.).

Доказательство.

Для борелевских множеств $B \in \mathcal{B}(\mathbb{R})$ положим

$$Q(B) = E(\xi \cdot I\{\eta \in B\}).$$

Заметим, что это заряд на вероятностном пространстве $(\mathbb{R}, \mathcal{B}(\mathbb{R}), P_\eta)$, абсолютно непрерывный относительно P_η . По теореме Радона-Никодима существует единственная P_η -п.н. случайная величина φ на $(\mathbb{R}, \mathcal{B}(\mathbb{R}), P_\eta)$ (борелевская функция!) такая, что для любого $B \in \mathcal{B}(\mathbb{R})$

$$Q(B) = \int_B \varphi(y) P_\eta(dy).$$

По определению получаем, что $\varphi(y) = E(\xi|\eta = y)$. ■

Вопрос: как связаны $E(\xi|\eta)$ и $E(\xi|\eta = y)$?

Утверждение

$$E(\xi|\eta = y) = \varphi(y) \iff E(\xi|\eta) = \varphi(\eta).$$

Доказательство.

\Rightarrow Пусть $B \in \mathcal{B}(\mathbb{R})$. Тогда

$$E(\xi \cdot I\{\eta \in B\}) = \int_B \varphi(y) P_\eta(dy) = E(\varphi(\eta) \cdot I\{\eta \in B\}).$$

Значит, $\varphi(\eta) = E(\xi|\eta)$.

\Leftarrow То же самое в обратном порядке:

$$E(\xi \cdot I\{\eta \in B\}) = E(\varphi(\eta) \cdot I\{\eta \in B\}) = \int_B \varphi(y) P_\eta(dy).$$

Следовательно, $\varphi(y) = E(\xi|\eta = y)$. ■

Свойства УМО

Условное математическое ожидание $E(\xi|\eta = y)$ обладает многими свойствами условного математического ожидания $E(\xi|\eta)$. В частности:

- линейность,
- сохранение отношения порядка,
- предельный переход под знаком УМО,
- неравенство Йенсена.

Доказываются они точно так же, как и свойства $E(\xi|\eta)$. В связи с этим детали мы опустим.

Условное распределение и условная плотность

Условные математические ожидания можно считать с помощью условных распределений. Правильное определение (регулярного) условного распределения - следующее.

Определение 10.3. Условным распределением случайной величины ξ относительно случайной величины η называется такая функция $P(B, y)$, $B \in \mathcal{B}(\mathbb{R})$, $y \in \mathbb{R}$, что

- 1) $P(\cdot, y)$ является распределением вероятностей на $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ для каждого $y \in \mathbb{R}$;
- 2) $P(B, \cdot)$ является борелевской функцией для каждого $B \in \mathcal{B}(\mathbb{R})$;
- 3) Для любых $B \in \mathcal{B}(\mathbb{R})$ и $A \in \mathcal{B}(\mathbb{R})$ выполняется равенство

$$P(\xi \in B, \eta \in A) = E(P(B, \eta) \cdot I\{\eta \in A\}).$$

Теорема 10.1. Условное распределение существует.

Обозначение: $P(B, y) = P(\xi \in B | \eta = y)$.

Замечание

Для любого $B \in \mathcal{B}(\mathbb{R})$ можно было бы определить

$$P(\xi \in B | \eta = y) := E(I\{\xi \in B\} | \eta = y).$$

Тогда будут выполнены свойства 2 и 3, но свойство 1 (грубо говоря, счетная аддитивность по B) не будет выполнено для всех множеств B сразу.

Определение 10.4. Условной плотностью случайной величины ξ относительно случайной величины η называется плотность условного распределения (по некоторой мере μ), то есть такая функция $f_{\xi|\eta}(x|y)$, $x, y \in \mathbb{R}$, что для любого $B \in \mathcal{B}(\mathbb{R})$ выполняется равенство:

$$P(\xi \in B | \eta = y) = \int_B f_{\xi|\eta}(x|y) \mu(dx).$$

Теорема о вычислении УМО

С помощью условных плотностей можно находить условные математические ожидания по той же формуле, что и обычные математические ожидания с помощью обычных плотностей.

Теорема 10.2. Если существует условная плотность $f_{\xi|\eta}(x|y)$ случайной величины ξ относительно случайной величины η , то для любой борелевской функции $g(x) : \mathbb{R} \rightarrow \mathbb{R}$ выполнено:

$$E(g(\xi)|\eta = y) = \int_{\mathbb{R}} g(x) f_{\xi|\eta}(x|y) \mu(dx).$$

Доказательство.

Пусть сначала $g(x) = I\{x \in A\}$ для некоторого борелевского множества A . Тогда для любого y

$$\begin{aligned} E(g(\xi)|\eta = y) &= P(\xi \in A|\eta = y) = \int_A f_{\xi|\eta}(x|y) \mu(dx) = \\ &= \int_{\mathbb{R}} I\{x \in A\} f_{\xi|\eta}(x|y) \mu(dx) = \int_{\mathbb{R}} g(x) f_{\xi|\eta}(x|y) \mu(dx). \end{aligned}$$

Раз формула выполнена для индикаторов, то в силу линейности математического ожидания формула верна и для линейных комбинация индикаторов, то есть для простых функций $g(x)$. Для произвольной функции $g(x)$ справедливость формулы устанавливается предельным переходом от простых функций. ■

Вычисление условной плотности

Вопрос: как вычислять условную плотность?

Теорема 10.3. (достаточное условие существования условной плотности)

Пусть случайные величины ξ и η имеют совместную плотность $f_{\xi,\eta}(x, y)$ по мере $\mu \times \nu$ на \mathbb{R}^2 . Тогда функция

$$f_{\xi|\eta}(x|y) = \begin{cases} \frac{f_{\xi,\eta}(x, y)}{f_{\eta}(y)}, & \text{если } f_{\eta}(y) > 0; \\ 0, & \text{иначе;} \end{cases}$$

является условной плотностью ξ по мере μ . Здесь $f_{\eta}(y)$ - это плотность случайной величины η по мере ν .

Доказательство.

Надо проверить, что функция $P(B, y) = \int_B f_{\xi|\eta}(x|y) \mu(dx)$ является условным распределением ξ относительно η , то есть проверить свойства 1-3 из определения.

Свойство 2 вытекает из теоремы Фубини, так как мы интегрируем совместную

плотность по одной переменной, а поскольку совместная плотность интегрируема на квадрате, то по теореме Фубини интеграл по одной координате является измеримой функцией (в нашем случае борелевской) от другой переменной.

Свойство 1 следует из свойств интеграла Лебега.

Проверим свойство 3. Для любых борелевских множеств A, B выполнено

$$\begin{aligned} P(\xi \in B, \eta \in A) &= \int_{B \times A} f_{\xi, \eta}(x, y) \mu(dx) \nu(dy) = \\ &= \int_A \left(\int_B \frac{f_{\xi, \eta}(x, y)}{f_{\eta}(y)} \mu(dx) \right) f_{\eta}(y) \nu(dy) = \int_A P(B, y) f_{\eta}(y) \nu(dy) = \\ &= \int_{\mathbb{R}} P(B, y) \cdot I\{\eta \in A\} f_{\eta}(y) \nu(dy) = E(P(B, \eta) \cdot I\{\eta \in A\}). \end{aligned}$$

Последнее равенство означает, что функция $\frac{f_{\xi, \eta}(x, y)}{f_{\eta}(y)}$ является условной плотностью. ■

Схема вычисления УМО

Зафиксируем схему вычисления УМО, которая у нас получилась. Необходимо вычислить $E(g(\xi)|\eta)$.

- 1) Находим $f_{\xi, \eta}(x, y)$ - совместную плотность ξ и η .
- 2) С ее помощью находим $f_{\xi|\eta}(x|y)$ - условную плотность ξ относительно η .
- 3) Вычисляем

$$h(y) = E(g(\xi)|\eta = y) = \int_{\mathbb{R}} g(x) f_{\xi|\eta}(x|y) dx.$$

- 4) Подставляем $E(g(\xi)|\eta) = h(\eta)$.

Пример

Случайные величины X и Y независимые $\text{Exp}(1)$. Найдите

$$E(X^2|X + Y) = ?$$

Решение

Сначала найдем совместную плотность $(X, X + Y)$. Для $B \in \mathcal{B}(\mathbb{R}^2)$ рассмотрим

$$\begin{aligned} P((X, X + Y) \in B) &= \int_{(x,y):(x,x+y) \in B} f_{X,Y}(x,y) dx dy = \\ &= \int_{(x,y):(x,x+y) \in B} e^{-x} e^{-y} \cdot I\{x, y > 0\} dx dy \quad \square \end{aligned}$$

сделаем замену $\alpha = x$, $\beta = x + y$, якобиан замены $= 1$

$$\square \int_{(\alpha,\beta) \in B} e^{-\beta} \cdot I\{\beta > \alpha > 0\} d\alpha d\beta.$$

Стало быть, функция $p_{X,X+Y}(\alpha, \beta) = e^{-\beta} \cdot I\{\beta > \alpha > 0\}$ и есть совместная плотность $(X, X + Y)$.

Теперь найдем плотность $X + Y$. Для этого совместную плотность $p_{X,X+Y}(\alpha, \beta)$ проинтегрируем по первой переменной.

$$p_{X+Y}(\beta) = \int_{\mathbb{R}} p(\alpha, \beta) d\alpha = \int_{\mathbb{R}} e^{-\beta} I\{\beta > \alpha > 0\} d\alpha = e^{-\beta} I\{\beta > 0\} \int_0^{\beta} d\alpha = \beta e^{-\beta} I\{\beta > 0\}.$$

Отсюда условная плотность равна

$$f_{X|X+Y}(\alpha|\beta) = \frac{p_{X,X+Y}(\alpha, \beta)}{p_{X+Y}(\beta)} = \frac{1}{\beta} I\{\beta > \alpha > 0\}.$$

Остается найти УМО

$$E(X^2|X + Y = \beta) = \int_{\mathbb{R}} \alpha^2 \cdot \frac{1}{\beta} I\{\beta > \alpha > 0\} d\alpha = \frac{\beta^2}{3}.$$

$$\text{В итоге, } E(X^2|X + Y) = \frac{(X + Y)^2}{3}.$$

Замечания

- 1) Если ξ и η независимы и имеют плотности, то условная плотность $f_{\xi|\eta}(x|y)$ равна просто плотности случайной величины ξ .
- 2) Вычисление $E(\xi|\eta)$ в случае векторного условия η происходит по тем же формулам.

- 3) Если условной плотности нет, то формально условное математическое ожидание вычисляется как интеграл по условному распределению

$$E(g(\xi)|\eta = y) = \int_{\mathbb{R}} g(x)Q_y(dx),$$

где $Q_y(B) = P(\xi \in B|\eta = y)$.

- 4) Особо важен случай вычисления условных математических ожиданий вида $E(g(\xi, \eta)|\eta)$ для независимых случайных величин ξ и η .

Важный случай

Лемма

Если ξ и η - независимые случайные величины, то

$$E(g(\xi, \eta)|\eta = y) = Eg(\xi, y).$$

Доказательство.

По теореме Фубини правая часть есть борелевская функция от y . Проверим интегральное свойство: для любого $B \in \mathcal{B}(\mathbb{R})$

$$E(g(\xi, \eta) \cdot I\{\eta \in B\}) = \int_{\mathbb{R}^2} g(x, y) \cdot I\{y \in B\} P_{\xi, \eta}(dx, dy) \equiv$$

пользуемся независимостью (совместное распределение - прямое произведение распределений)

$$\equiv \int_B \left(\int_{\mathbb{R}} g(x, y) P_{\xi}(dx) \right) P_{\eta}(dy) = \int_B Eg(\xi, y) P_{\eta}(dy).$$

■

Теорема о наилучшем квадратичном прогнозе

Теорема 10.4. (о наилучшем квадратичном прогнозе)

Пусть ξ - случайная величина на (Ω, \mathcal{F}, P) с конечным вторым моментом, \mathcal{C} - под- σ -алгебра в \mathcal{F} ю Обозначим через $\mathcal{A}_{\mathcal{C}}$ - множество всех \mathcal{C} -измеримых случайных величин. Тогда

$$E(\xi - E(\xi|\mathcal{C}))^2 = \min_{\eta \in \mathcal{A}_{\mathcal{C}}} E(\xi - \eta)^2,$$

причем $E(\xi|\mathcal{C})$ - это единственный минимум.

Доказательство.

Пусть $\eta \in \mathcal{A}_{\mathfrak{C}}$. Тогда

$$\begin{aligned} E(\xi - \eta)^2 &= E(\xi - E(\xi|\mathfrak{C}) + E(\xi|\mathfrak{C}) - \eta)^2 = \\ &= E(\xi - E(\xi|\mathfrak{C}))^2 + E(E(\xi|\mathfrak{C}) - \eta)^2 + 2E\left((\xi - E(\xi|\mathfrak{C}))(E(\xi|\mathfrak{C}) - \eta)\right). \end{aligned}$$

Проверим, что третье слагаемое равно нулю. Тогда $E(\xi - \eta)^2$ минимизируется, когда второй квадрат равен нулю (т.к. только он зависит от η). Действительно, по формуле полной вероятности

$$E\left((\xi - E(\xi|\mathfrak{C}))(E(\xi|\mathfrak{C}) - \eta)\right) = E\left(E((\xi - E(\xi|\mathfrak{C}))(E(\xi|\mathfrak{C}) - \eta)|\mathfrak{C})\right) \equiv$$

свойство (9) УМО

$$\equiv E\left((E(\xi|\mathfrak{C}) - \eta)E(\xi - E(\xi|\mathfrak{C})|\mathfrak{C})\right) = 0.$$

В итоге,

$$E(\xi - \eta)^2 = E(\xi - E(\xi|\mathfrak{C}))^2 + E(E(\xi|\mathfrak{C}) - \eta)^2 \geq E(\xi - E(\xi|\mathfrak{C}))^2,$$

причем равенство достигается тогда и только тогда, когда $\eta = E(\xi|\mathfrak{C})$. ■

11.

Лекция 11

Постановка задачи

Напомним байесовский подход к сравнению оценок.

- Пусть \mathbf{X} - наблюдение с неизвестным распределением P , принадлежащем параметрическому семейству $\{P_\theta, \theta \in \Theta\}$.
- Пусть задано некоторое распределение вероятностей Q (которое называется априорным) на множестве значений параметра Θ .
- Пусть задана функция потерь $\rho(x, y), x, y \in \Theta$.
- Для каждой оценки $\theta^*(\mathbf{X})$ определим функцию риска

$$R_{\theta^*}(\theta) = E_\theta \rho(\theta^*(\mathbf{X}), \theta).$$

Задача Требуется найти такую оценку $\hat{\theta}(\mathbf{X})$, что

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

Предположения

Оказывается, что при некоторых условиях можно получить явное выражение для наилучшей оценки. Для этого факта сделаем несколько предположений.

- 1) Будем считать, что параметр числовой: $\Theta \in \mathbb{R}$.
- 2) Распределение Q имеет плотность $q(t)$ по мере λ .
- 3) Функция потерь является квадратичной и, значит,

$$R_{\theta^*}(\theta) = E_\theta (\theta^*(\mathbf{X}) - \theta)^2$$

- 4) Семейство распределений $\{P_\theta, \theta \in \Theta\}$ — это доминируемое семейство распределений с плотностью $p_\theta(X)$ по мере μ .

В рамках этих предположений:

Определение 11.1. Плотность $q(t)$ называется априорной плотностью параметра θ .

Байесовские оценки

Определение 11.2. Функция

$$q(t|\mathbf{X}) = \frac{q(t)p_t(\mathbf{X})}{\int_{\Theta} q(s)p_s(\mathbf{X})\lambda(ds)}$$

называется апостериорной плотностью параметра θ (плотностью при заданном \mathbf{X}).

Определение 11.3. Оценка

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

называется байесовской оценкой параметра θ для заданного априорного распределения Q .

Мы готовы сформулировать теорему о байесовской оценке, которая говорит, что $\hat{\theta}_Q(\mathbf{X})$ это ровно то, что мы ищем.

Теорема 11.1. (о байесовской оценке)

В приведенных условиях байесовская оценка $\hat{\theta}_Q(\mathbf{X})$ является наилучшей оценкой параметра θ в байесовском подходе.

Доказательство.

Идея доказательства:

Доказательство будет сводиться к теореме о наилучшем квадратичном прогнозе. Мы хотим показать, что байесовская оценка - это будет условное математическое ожидание в некотором пространстве и, соответственно, она будет величиной, которая минимизирует квадратичный прогноз. В байесовском подходе мы рассматриваем параметр как случайную величину - у него есть распределение вероятности Q . И идея состоит в том, чтобы рассмотреть пару: параметр и наблюдение как случайный вектор и, соответственно, понять, какую они имеют совместную плотность. Далее, показать, что интегральное значение функции риска для оценки совпадет с математическим ожиданием от квадрата разности "оценка минус параметр но уже в некотором другом вероятностном пространстве, в котором параметр будет случайной величиной.

Итак, рассмотрим на пространстве $\Theta \times \mathcal{X}$ следующую функцию:

$$f(t, \mathbf{x}) = q(t)p_t(\mathbf{x}), t \in \Theta, \mathbf{x} \in \mathcal{X}.$$

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 $f(t, \mathbf{x})$ - это плотность некоторого вероятностного распределения на $\Theta \times \mathcal{X}$ по мере $\lambda \times \mu$. Обозначим через \tilde{P} соответствующую вероятностную меру.

Далее, посмотрим на вектор (θ, \mathbf{X}) , как на случайный вектор на вероятностном пространстве $(\Theta \times \mathcal{X}, \tilde{P})$, заданный по правилу:

$$(\theta, \mathbf{X})(t, \mathbf{x}) = (t, \mathbf{x}).$$

Мы хотим минимизировать по θ^* следующий функционал:

$$\begin{aligned} \int_{\Theta} R_{\theta^*}(t) q(t) \lambda(dt) &= \int_{\Theta} E_t(\theta^*(\mathbf{X}) - t)^2 q(t) \lambda(dt) = \\ &= \int_{\Theta} \left(\int_{\mathcal{X}} (\theta^*(\mathbf{x}) - t)^2 p_t(\mathbf{x}) \mu(d\mathbf{x}) \right) q(t) \lambda(dt) = \\ &= \int_{\Theta \times \mathcal{X}} (\theta^*(\mathbf{x}) - t)^2 f(t, \mathbf{x}) \lambda(dt) \times \mu(d\mathbf{x}) = \tilde{E}(\theta^*(\mathbf{X}) - \theta)^2. \end{aligned}$$

Тем самым, задача поиска наилучшей оценки в байесовском подходе свелась к минимизации математического ожидания в пространстве \tilde{E} :

$$\tilde{E}(\theta^*(\mathbf{X}) - \theta)^2 \rightarrow \min_{\theta^*(\mathbf{X})}.$$

Теорема о наилучшем квадратичном прогнозе говорит, что решением такой задачи является $\tilde{E}(\theta|\mathbf{X})$. Теперь осталось показать, что $\hat{\theta}_Q(\mathbf{X}) = \tilde{E}(\theta|\mathbf{X})$.

Для этого заметим, что:

- $q(t)$ есть плотность θ ;
- Наблюдение \mathbf{X} имеет плотность

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

- Условная плотность \mathbf{X} относительно θ равна $\frac{f(t, \mathbf{x})}{q(t)} = p_t(\mathbf{x})$.

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

$$\frac{f(t, \mathbf{x})}{g(\mathbf{x})} = \frac{q(t) p_t(\mathbf{x})}{\int_{\Theta} q(s) p_s(\mathbf{x}) \lambda(ds)} = q(t|\mathbf{x}),$$

т.е. это и есть апостериорная плотность.

Применяя формулу для вычисления УМО, получаем, что

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

Значит, $\hat{\theta}_Q(\mathbf{x}) = \tilde{E}(\theta|\mathbf{X})$. ■

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из нормального распределения $\mathcal{N}(\theta, 1)$, $\theta \in \mathbb{R}$. Найдите байесовскую оценку параметра θ , если априорное распределение θ есть $\mathcal{N}(a, \sigma^2)$.

Решение

У нас есть формулы, надо просто посчитать. Найдем плотность выборки:

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

Теперь рассмотрим произведение $p_t(\mathbf{X})q(t)$:

$$p_t(\mathbf{X})q(t) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left(-\frac{1}{2} \sum_{i=1}^n X_i^2 + t \sum_{i=1}^n X_i - \frac{n}{2}t^2\right) \times \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(t-a)^2}{2\sigma^2}\right).$$

Как функция от t это произведение представляет собой ненормированную плотность нормального распределения.

Найдем параметры этого распределения. Если убрать все слагаемые, не зависящие от t в экспоненте, то останется:

$$-\left(\frac{n}{2} + \frac{1}{2\sigma^2}\right)t^2 + t\left(\sum_{i=1}^n X_i + \frac{a}{\sigma^2}\right) = -\frac{1}{2}\left(n + \frac{1}{\sigma^2}\right)\left(t - \frac{\bar{\mathbf{X}} + \frac{a}{n\sigma^2}}{1 + \frac{1}{n\sigma^2}}\right)^2 + \dots$$

Таким образом, апостериорное распределение есть $\mathcal{N}\left(\frac{\bar{\mathbf{X}} + \frac{a}{n\sigma^2}}{1 + \frac{1}{n\sigma^2}}, \left(n + \frac{1}{\sigma^2}\right)^{-1}\right)$.

Байесовская оценка - это математическое ожидание апостериорного распределения, стало быть, она равна:

$$\hat{\theta}_Q(\mathbf{X}) = \frac{\bar{\mathbf{X}} + \frac{a}{n\sigma^2}}{1 + \frac{1}{n\sigma^2}}.$$

Стоит отметить, что с ростом n байесовская оценка становится похожей на $\bar{\mathbf{X}}$

Минимаксный подход

Напомним про минимаксный подход к сравнению оценок. В нём задача поиска наилучшей оценки ставится следующим образом:

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

т.е. мы хотим минимизировать максимальное значение функции риска.

Оказывается, минимаксный подход тесно связан с байесовским, а минимаксная оценка в некоторых случаях может быть вычислена как байесовская для правильно подобранного априорного распределения Q .

Начнем обсуждение со следующего достаточного условия минимаксности оценки.

Достаточное условие минимаксности оценки

Лемма (1) Пусть $\hat{\theta}(\mathbf{X})$ - такая оценка параметра θ , что существует вероятностное распределение Q на Θ со следующим условием: для всех $\theta \in \Theta$

$$R_{\hat{\theta}}(\theta) \leq \int_{\Theta} R_{\hat{\theta}_Q}(t) Q(dt),$$

где $\hat{\theta}_Q(\mathbf{X})$ - это байесовская оценка θ для априорного распределения Q .

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

Доказательство.

Пусть $\theta^*(\mathbf{X})$ - любая оценка θ . Тогда

$$\sup_{\theta \in \Theta} R_{\theta^*}(\theta) \geq \int_{\Theta} R_{\theta^*}(t) Q(dt) \geq \int_{\Theta} R_{\hat{\theta}_Q}(t) Q(dt) \geq \sup_{\theta \in \Theta} R_{\hat{\theta}}(\theta).$$

■

Замечание

Отметим, что в условиях леммы 1 должно выполняться равенство

$$R_{\hat{\theta}}(\theta) = \int_{\Theta} R_{\hat{\theta}_Q}(t) Q(dt)$$

почти наверное по мере Q .

Следствие (1)

Если оценка $\hat{\theta}(\mathbf{X})$ удовлетворяет двум условиям:

- $\hat{\theta}(\mathbf{X})$ является байесовской оценкой θ для некоторого априорного распределения Q ;
- существует такая константа c , что функция риска $R_{\hat{\theta}}(\theta) = c Q$ - почти наверное и $R_{\hat{\theta}}(\theta) \leq c$ для всех θ ;

то оценка $\hat{\theta}(\mathbf{X})$ является минимаксной.

Выводы

Первый вывод.

Минимаксная оценка - это байесовская оценка, которая "выравнивает" функцию риска.

Вопрос: но какое априорное распределение надо взять?

Второй вывод.

Если оценка $\hat{\theta}(\mathbf{X}) = \hat{\theta}_Q(\mathbf{X})$ удовлетворяет условию следствия, то для любого априорного распределения \tilde{Q} будет выполнено следующее неравенство:

$$c = \int_{\Theta} R_{\hat{\theta}}(t) Q(dt) \geq \int_{\Theta} R_{\hat{\theta}}(t) \tilde{Q}(dt).$$

Но

$$\int_{\Theta} R_{\hat{\theta}}(t) \tilde{Q}(dt) \geq \int_{\Theta} R_{\hat{\theta}_{\tilde{Q}}}(t) \tilde{Q}(dt), \quad \int_{\Theta} R_{\hat{\theta}}(t) \tilde{Q}(dt) \leq \int_{\Theta} R_{\hat{\theta}_Q}(t) Q(dt).$$

Таким образом, Q - это то распределение, для которого среднее значение функции риска соответствующей байесовской оценки является максимальным.

$$\int_{\Theta} R_{\hat{\theta}_Q}(t) Q(dt) = \max_{\tilde{Q}} \int_{\Theta} R_{\hat{\theta}_{\tilde{Q}}}(t) \tilde{Q}(dt).$$

Определение 11.4. Подобное априорное распределение называется наихудшим априорным распределением параметра θ .

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из распределения Бернулли $\text{Bin}(1, \theta)$, $\theta \in [0, 1]$. Найдите минимаксную оценку θ .

Решение

Наш план состоит в следующем:

- найти оценку с постоянной функцией риска,

- доказать, что подобная оценка является байесовской для некоторого априорного распределения.

Разумно искать оценку в виде $\hat{\theta}(\mathbf{X}) = a\bar{\mathbf{X}} + b$ для $a, b \in \mathbb{R}$. Несложными вычислениями можно получить, что

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

Но как искать теперь наилучшее априорное распределение?

Двойственные распределения

Идея: надо посмотреть на функцию правдоподобия как на (ненормированную) плотность от параметра θ .

В нашем примере

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

По $\theta \in [0, 1]$ - это ненормированная плотность бета-распределения. В связи с этим можно взять в качестве априорного распределения тоже бета-распределение.

Возьмем априорное распределение $\text{Beta}(\alpha, \alpha)$. Тогда априорная плотность будет равна $q(t) = \frac{t^{\alpha-1}(1-t)^{\alpha-1}}{B(\alpha, \alpha)}$. Отсюда

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

Следовательно, априорное распределение тоже есть бета-распределение $\text{Beta}(n\bar{\mathbf{X}} + \alpha, n(1 - \bar{\mathbf{X}}) + \alpha)$.

Пример

Следовательно, байесовская оценка будет иметь вид:

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

Наша оценка имеет:

$$\hat{\theta}(\mathbf{X}) = a\bar{\mathbf{X}} + b = \left(\frac{1}{1 + \frac{1}{\sqrt{n}}} \right) \bar{\mathbf{X}} + \frac{1}{2(\sqrt{n} + 1)} = \frac{\bar{\mathbf{X}} + \frac{1}{2\sqrt{n}}}{1 + \frac{1}{\sqrt{n}}}.$$

Остается заметить, что при $\alpha = \frac{\sqrt{n}}{2}$ оценки совпадают, и, значит, $\hat{\theta}(\mathbf{X})$ есть байесовская оценка для некоторого априорного распределения.

Еще одно достаточное условие

Однако, наихудшее априорное распределение есть не всегда! В этом случае может помочь второе достаточное условие того, что оценка минимаксная.

Лемма (2)

Пусть $\hat{\theta}(\mathbf{X})$ - такая оценка параметра θ , что существует последовательность априорных распределений $\{Q_k, k \in \mathbb{N}\}$ на Θ со следующим условием: для всех $\theta \in \Theta$

$$R_{\hat{\theta}}(\theta) \leq \overline{\lim}_{k \rightarrow +\infty} \int_{\Theta} R_{\hat{\theta}_{Q_k}}(t) Q_k(dt).$$

Тогда оценка $\hat{\theta}(\mathbf{X})$ является минимаксной.

Доказательство. Пусть $\theta^*(\mathbf{X})$ - любая оценка θ . Тогда для любого $k \in \mathbb{N}$

$$\sup_{\theta \in \Theta} R_{\theta^*}(\theta) \geq \int_{\Theta} R_{\theta^*}(t) Q_k(dt) \geq \int_{\Theta} R_{\hat{\theta}_{Q_k}}(t) Q_k(dt).$$

Перейдем к пределу по k . Тогда

$$\sup_{\theta \in \Theta} R_{\theta^*}(\theta) \geq \overline{\lim}_{k \rightarrow +\infty} \int_{\Theta} R_{\hat{\theta}_{Q_k}}(t) Q_k(dt) \geq \sup_{\theta \in \Theta} R_{\hat{\theta}}(\theta),$$

что и доказывает минимаксность $\hat{\theta}(\mathbf{X})$. ■

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из нормального распределения $\mathcal{N}(\theta, 1)$, $\theta \in \mathbb{R}$. Найдите минимаксную оценку параметра θ .

Решение

План решения. Заметим, что \bar{X} проходит в качестве претендента на минимаксность, так как его функция риска не зависит от параметра.

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

Однако, сам \bar{X} является байесовской оценкой ни для какого априорного распределения, так что мы хотим подобрать подходящую последовательность априорных распределений с условием:

$$\overline{\lim}_{k \rightarrow +\infty} \int_{\Theta} R_{\hat{\theta}_{Q_k}}(t) Q_k(dt) \geq \frac{1}{n}.$$

Возьмем $Q_k = \mathcal{N}(0, k)$. Тогда мы знаем, что

$$\hat{\theta}_{Q_k}(\mathbf{X}) = \frac{\bar{\mathbf{X}}}{1 + \frac{1}{nk}}$$

Найдем искомый интеграл от функции риска:

$$\begin{aligned} \int_Q R_{\hat{\theta}_{Q_k}}(t) Q_k(dt) &= \int_{\mathbb{R}} E_t \left(\frac{\bar{\mathbf{X}}}{1 + \frac{1}{nk}} - t \right)^2 q_k(t) dt = \\ &= \int_{\mathbb{R}} \left[\int_{\mathbb{R}^n} \left(\frac{\bar{\mathbf{X}}}{1 + \frac{1}{nk}} - t \right)^2 p_t(x_1, \dots, x_n) dx_1 \dots dx_n \right] q_k(t) dt = \end{aligned}$$

(выразим $q_k(t|\mathbf{x}) = \frac{q_k(t)p_t(\mathbf{x})}{\tilde{p}(\mathbf{x})}$, где $\tilde{p}(\mathbf{x}) = \int_{\mathbb{R}} q_k(t)p_t(\mathbf{x})dt$)

(переставим интегралы)

$$= \int_{\mathbb{R}^n} \left[\int_{\mathbb{R}} \left(\frac{\bar{\mathbf{X}}}{1 + \frac{1}{nk}} - t \right)^2 q_k(t|\mathbf{x}) dt \right] \tilde{p}(x_1, \dots, x_n) dx_1 \dots dx_n =$$

(внутренний интеграл есть дисперсия апостериорного распределения)

$$= \frac{1}{n + \frac{1}{k}} \rightarrow \frac{1}{n} \text{ при } k \rightarrow +\infty.$$

Стало быть, $\bar{\mathbf{X}}$ - это минимаксная оценка θ .

12.

Лекция 12

Понятие оптимальной оценки

Ранее мы подробно рассмотрели асимптотический, байесовский и минимаксный подходы к сравнению оценок. Теперь осталось рассмотреть самый "сильный" подход - равномерный. Как и в случае байесовского и минимаксного подходов будем изучать только квадратичную функцию потерь. Напомним, что в среднеквадратичном подходе 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. Будем работать в классе несмещённых оценок. В таком случае $E_{\theta}(\hat{\theta}(\mathbf{X}) - \theta)^2 = D_{\theta}\hat{\theta}(\mathbf{X})$. Более того, можно оценивать не только сам θ , но и функции от него.

Все это приводит нас к определению оптимальной оценки

Определение 12.1. Пусть \mathbf{X} - наблюдение с неизвестным распределением $P \in \{P_{\theta}, \theta \in \Theta\}$, а $\tau(\theta) \in \mathbb{R}^k$ - некоторая функция от θ . Тогда оценка $\hat{\tau}(\mathbf{X})$ называется оптимальной оценкой $\tau(\theta)$, если она является несмещённой оценкой $\tau(\theta)$ и имеет равномерно наименьшую дисперсию, т.е. для любой другой несмещённой оценки $\theta^*(\mathbf{X})$ выполнено

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

В многомерном случае, т.е. при $\hat{\tau}(\mathbf{X}) \in \mathbb{R}^k, k \geq 2$, неравенство дисперсий означает, что разность матриц ковариаций является неотрицательно определенной матрицей:

$$D_{\theta}\hat{\tau}(\mathbf{X}) \leq D_{\theta}\theta^*(\mathbf{X}) \iff D_{\theta}\theta^*(\mathbf{X}) - D_{\theta}\hat{\tau}(\mathbf{X}) \text{ неотрицательно определена.}$$

Замечание

Эффективная оценка $\tau(\theta)$ является оптимальной, но она существует далеко не для всех функций от θ .

Достаточные статистики

Определение 12.2. Пусть \mathbf{X} - наблюдение с неизвестным распределением $P \in \{P_{\theta}, \theta \in \Theta\}$. Статистика $S(\mathbf{X})$ называется достаточной для семейства $\{P_{\theta}, \theta \in \Theta\}$,

если существует вариант условного распределения $P_\theta(\mathbf{X} \in B | S(\mathbf{X}) = s)$, который не зависит от параметра θ , т.е. существует функция $P(B, s)$ такая, что для всех $\theta \in \Theta$

$$P_\theta(\mathbf{X} \in B | S(\mathbf{X}) = s) = P(B, s) P_\theta^S - \text{почти наверное,}$$

где P_θ^S - это распределение статистики $S(\mathbf{X})$, когда \mathbf{X} имеет распределение P_θ .

Наблюдение: сам \mathbf{X} является достаточной статистикой, но интересен случай, когда размерность $S(\mathbf{X})$ совпадает с размерностью параметра.

Ответ на вопрос, зачем нужны достаточные статистики, дает следующая теорема.

Теорема 12.1. Колмогорова-Блэкуэлла-Рао

Пусть $\hat{\theta}(\mathbf{X})$ - несмещённая оценка $\tau(\theta) \in \mathbb{R}$ с конечным вторым моментом: $E_\theta(\hat{\theta}(\mathbf{X}))^2 < +\infty$ для всех $\theta \in \Theta$. Далее, пусть $S(\mathbf{X})$ - достаточная статистика для семейства распределений $\{P_\theta, \theta \in \Theta\}$. Тогда

- 1) $\theta^*(\mathbf{X}) = E_\theta(\hat{\theta}(\mathbf{X}) | S(\mathbf{X}))$ есть несмещённая оценка $\tau(\theta)$;
- 2) для всех $\theta \in \Theta$ выполнено $D_\theta \theta^*(\mathbf{X}) \leq D_\theta \hat{\theta}(\mathbf{X})$;
- 3) равенство в неравенстве выше для всех $\theta \in \Theta$ достигается тогда и только тогда, когда $\hat{\theta}(\mathbf{X})$ является $S(\mathbf{X})$ -измеримой (т.е. измеримой функцией от $S(\mathbf{X})$).

Тем самым, теорема позволяет "улучшить" оценку, получив другую с равномерно меньшей дисперсией.

Доказательство.

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

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

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

- 2) Теперь проверим соотношение между дисперсиями. Для этого заметим, что функция $h(x) = (x - \tau(\theta))^2$ является выпуклой книзу. Тогда можно воспользоваться неравенством Йенсена для УМО:

$$(\theta^*(\mathbf{X}) - \tau(\theta))^2 = (E_\theta(\hat{\theta}(\mathbf{X}) | S(\mathbf{X})) - \tau(\theta))^2 \leq E_\theta((\hat{\theta}(\mathbf{X}) - \tau(\theta))^2 | S(\mathbf{X})). \quad (1)$$

Возьмем математическое ожидание E_θ от обеих частей неравенства:

$$D_\theta \theta^*(\mathbf{X}) = E_\theta(\theta^*(\mathbf{X}) - \tau(\theta))^2 \leq E_\theta(\hat{\theta}(\mathbf{X}) - \tau(\theta))^2 = D_\theta \hat{\theta}(\mathbf{X}).$$

- 3) Осталось понять, когда будет выполняться равенство в данном неравенстве. Заметим, что разность правой и левой части в (1) равна:

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

В итоге, разность дисперсий

$$D_{\theta}\hat{\theta}(\mathbf{X}) - D_{\theta}\theta^*(\mathbf{X}) = E_{\theta}(E_{\theta}((\hat{\theta}(\mathbf{X}) - \theta^*(\mathbf{X}))^2 | S(\mathbf{X}))) = E_{\theta}(\hat{\theta}(\mathbf{X}) - \theta^*(\mathbf{X}))^2,$$

равна нулю тогда и только тогда, когда $\hat{\theta}(\mathbf{X}) = \theta^*(\mathbf{X})$ P_{θ} - почти наверное для всех $\theta \in \Theta$. ■

Следствие (1)

Теорема Колмогорова-Блэкуэлла-Рао верная и в многомерном случае.

Доказательство.

Пусть $\tau(\theta) \in \mathbb{R}^k$. Возьмем произвольный ненулевой вектор $a \in \mathbb{R}^k$. Тогда $\langle \hat{\theta}(\mathbf{X}), a \rangle$ есть несмещённая оценка $\langle \tau(\theta), a \rangle$. Следовательно, в силу уже доказанного

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

Но это означает, что

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

Следовательно, $D_{\theta}\theta^*(\mathbf{X}) \leq D_{\theta}\hat{\theta}(\mathbf{X})$. Критерий равенства доказывается аналогично одномерному случаю. ■

Вопрос: при каких условиях оценка в теореме Колмогорова-Блэкуэлла-Рао является оптимальной?

Следствие (2)

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 оптимальная. Тогда найдется оценка $\hat{\theta}(\mathbf{X})$, которая будет лучше, т.е. дисперсия $\hat{\theta}(\mathbf{X})$ будет равномерно не больше дисперсии $\theta^*(\mathbf{X})$ и для некоторого $\theta_0 \in \Theta$

$$D_{\theta_0}\hat{\theta}(\mathbf{X}) < D_{\theta_0}\theta^*(\mathbf{X}).$$

Но тогда по теореме Колмогорова-Блэкуэлла-Рао оценка

$$\theta^{**}(\mathbf{X}) = E_{\theta}(\hat{\theta}(\mathbf{X}) | S(\mathbf{X}))$$

будет не хуже, чем $\hat{\theta}(\mathbf{X})$. Однако и $\theta^*(\mathbf{X})$, и $\theta^{**}(\mathbf{X})$ являются функциями от $S(\mathbf{X})$. Следовательно, они равны, и $\theta^*(\mathbf{X})$ не хуже, чем $\hat{\theta}(\mathbf{X})$. Противоречие. Значит, $\theta^*(\mathbf{X})$ оптимальна. ■

Полные статистики

Определение 12.3. 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 - почти наверное для всех $\theta \in \Theta$.

Смысл определения: несмещённо оценить нуль можно лишь нулевой функцией от $S(\mathbf{X})$.

Оказывается, что если статистика $S(\mathbf{X})$ является и достаточной, и полной, то улучшение оценки в теореме Колмогорова-Блэкуэлла-Рао сразу даст нам оптимальную оценку.

Теорема 12.2. (об оптимальной оценке)

Пусть $S(\mathbf{X})$ - полная достаточная статистика для семейства $\{P_\theta, \theta \in \Theta\}$. Пусть функция $\varphi(s)$ такова, что $\varphi(S(\mathbf{X}))$ - несмещённая оценка $\tau(\theta)$ и $E_\theta \varphi^2(S(\mathbf{X})) < +\infty$ для $\forall \theta \in \Theta$. Тогда $\varphi(S(\mathbf{X}))$ является оптимальной оценкой $\tau(\theta)$.

Доказательство. Согласно следствию 2 достаточно проверить, что $\varphi(S(\mathbf{X}))$ является единственной $S(\mathbf{X})$ -измеримой несмещённой оценкой $\tau(\theta)$. Пусть $\psi(S(\mathbf{X}))$ - другая несмещённая оценка $\tau(\theta)$. Тогда для любого $\theta \in \Theta$

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

В силу полноты статистики $S(\mathbf{X})$ получаем, что $\varphi(S(\mathbf{X})) = \psi(S(\mathbf{X}))$ P_θ - почти наверное. Отсюда получаем, что $\varphi(S(\mathbf{X}))$ - это единственная $S(\mathbf{X})$ -измеримая несмещённая оценка $\tau(\theta)$. Следовательно, она оптимальна. ■

Замечания

Замечание (1)

Если $S(\mathbf{X})$ - полная достаточная статистика, то для нахождения оптимальной оценки $\tau(\theta)$ достаточно решить (относительно φ) уравнение несмещённости

$$E_\theta \varphi(S(\mathbf{X})) = \tau(\theta).$$

Замечание (2)

Если $\tau(\theta) \in \mathbb{R}^k$ при $k > 1$, то оптимальная оценка вектора есть вектор из оптимальных оценок компонент.

Вопрос 1: как находить достаточные статистики?

Вопрос 2: как проверить их на полноту?

Критерий факторизации

На первый вопрос даёт ответ следующая теорема.

Теорема 12.3. (критерий факторизации Неймана-Фишера)

Пусть наблюдение \mathbf{X} имеет неизвестное распределение $P \in \{P_\theta, \theta \in \Theta\}$, и это семейство является доминируемым относительно меры μ с плотностью $p_\theta(\mathbf{x})$. Тогда статистика $S(\mathbf{X})$ будет достаточной для семейства $\{P_\theta, \theta \in \Theta\}$ тогда и только тогда, когда существует представление p_θ в виде

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

где ψ_θ и h - некоторые неотрицательные измеримые функции.

Таким образом, статистика является достаточной тогда и только тогда, когда плотность зависит от параметра через функцию от нее.

Доказательство.(для дискретного случая)

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

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

Первый множитель в произведении не зависит от θ в силу 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 $h(\mathbf{x}) = P_\theta(\mathbf{X} = \mathbf{x} | S(\mathbf{X}) = S(\mathbf{x}))$, а $\psi_\theta(S(\mathbf{x})) = P_\theta(S(\mathbf{X}) = S(\mathbf{x}))$.

\Leftarrow Теперь предположим, что существует представление $P_\theta(\mathbf{X} = \mathbf{x})$ в виде $h(\mathbf{x})\psi_\theta(S(\mathbf{x}))$. Рассмотрим условное распределение $P_\theta(\mathbf{X} = \mathbf{x} | S(\mathbf{X}) = s)$. Если $S(\mathbf{x}) \neq s$, то такая вероятность равна нулю. Иначе:

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

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

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

Формула пересчета УМО

Для доказательства общего случая нам понадобится формула пересчета УМО. Она аналогичная формуле пересчета обычных математических ожиданий (интегралов Лебега): если мера \tilde{P} абсолютно непрерывна относительно P , то

$$\tilde{E}\xi = E\left(\xi \cdot \frac{d\tilde{P}}{dP}\right).$$

Теорема 12.4. Пусть (Ω, \mathcal{F}, P) - вероятностное пространство, \tilde{P} - другая вероятностная мера, абсолютно непрерывная относительно P ($\tilde{P} \ll P$), а $\frac{d\tilde{P}}{dP}$ - соответствующая производная Радона-Никодима. Тогда для любой случайной величины ξ с условием $E|\xi| < +\infty$ и любой под- σ -алгебры $\mathfrak{C} \subset \mathcal{F}$ выполнено

$$\tilde{E}(\xi|\mathfrak{C}) = \frac{E\left(\xi \cdot \frac{d\tilde{P}}{dP}|\mathfrak{C}\right)}{E\left(\frac{d\tilde{P}}{dP}|\mathfrak{C}\right)}.$$

Доказательство. (формулы пересчета)

Заметим, что $E(\frac{d\tilde{P}}{dP}|\mathfrak{C}) \neq 0$ \tilde{P} - почти наверное. Обозначим $B = \{E(\frac{d\tilde{P}}{dP}|\mathfrak{C}) = 0\} \in \mathfrak{C}$.

Тогда, применяя интегральное свойство УМО, получаем:

$$0 = E\left(I_B \cdot E\left(\frac{d\tilde{P}}{dP}|\mathfrak{C}\right)\right) = E\left(I_B \cdot \frac{d\tilde{P}}{dP}\right) = \tilde{E}I_B = \tilde{P}(B).$$

Теперь докажем саму формулу. Правая часть является \mathfrak{C} -измеримой случайной величины, проверим интегральное свойство: для любого $A \in \mathfrak{C}$

$$\tilde{E}\left(I_A \cdot \frac{E(\xi \cdot \frac{d\tilde{P}}{dP}|\mathfrak{C})}{E(\frac{d\tilde{P}}{dP}|\mathfrak{C})}\right) = E\left(I_A \cdot \frac{E(\xi \cdot \frac{d\tilde{P}}{dP}|\mathfrak{C})}{E(\frac{d\tilde{P}}{dP}|\mathfrak{C})} \cdot \frac{d\tilde{P}}{dP}\right) =$$

(интегральное свойство)

$$= E\left(I_A \cdot E\left(\frac{E(\xi \cdot \frac{d\tilde{P}}{dP}|\mathfrak{C})}{E(\frac{d\tilde{P}}{dP}|\mathfrak{C})} \cdot \frac{d\tilde{P}}{dP} \middle| \mathfrak{C}\right)\right) =$$

(применяем свойство 9 УМО)

$$= E\left(I_A \cdot \frac{E(\xi \cdot \frac{d\tilde{P}}{dP} | \mathfrak{C})}{E(\frac{d\tilde{P}}{dP} | \mathfrak{C})} \cdot E\left(\frac{d\tilde{P}}{dP} | \mathfrak{C}\right)\right) = E\left(I_A \cdot E\left(\xi \cdot \frac{d\tilde{P}}{dP} | \mathfrak{C}\right)\right) =$$

(интегральное свойство)

$$= E\left(I_A \cdot \xi \cdot \frac{d\tilde{P}}{dP}\right) = \tilde{E}(I_A \cdot \xi).$$

Все доказано. ■

Доказательство. (критерия факторизации)

Зафиксируем некоторое $\theta' \in \Theta$ и для любого $\theta \neq \theta'$ введем вероятностную меру

$$P_{\theta, \theta'} = \frac{1}{2}P_{\theta} + \frac{1}{2}P_{\theta'}.$$

Заметим, что $P_{\theta}, P_{\theta'} \ll P_{\theta, \theta'}$, поэтому определены плотности:

$$f_{\theta}(\mathbf{X}) = \frac{dP_{\theta}}{dP_{\theta, \theta'}} = \frac{p_{\theta}(\mathbf{X})}{\frac{1}{2}(p_{\theta}(\mathbf{X}) + p_{\theta'}(\mathbf{X}))},$$

$$f'_{\theta'}(\mathbf{X}) = \frac{dP_{\theta'}}{dP_{\theta, \theta'}} = \frac{p_{\theta'}(\mathbf{X})}{\frac{1}{2}(p_{\theta}(\mathbf{X}) + p_{\theta'}(\mathbf{X}))}.$$

\Rightarrow Пусть $S(\mathbf{X})$ - достаточная статистика. Покажем, что $f_{\theta}(\mathbf{X})$ есть функция от $S(\mathbf{X})$, $f_{\theta}(\mathbf{X}) = \psi_{\theta}(S(\mathbf{X}))$. Этого будет достаточно для доказательства существования нужного представления для $p_{\theta}(\mathbf{x})$:

$$p_{\theta}(\mathbf{x}) = \frac{f_{\theta}(\mathbf{x})p_{\theta'}(\mathbf{x})}{2 - f_{\theta}(\mathbf{x})} = \frac{\psi_{\theta}(S(\mathbf{x}))}{2 - \psi_{\theta}(S(\mathbf{x}))} \cdot p_{\theta'}(\mathbf{x}),$$

где $h(x) = p_{\theta'}(\mathbf{x})$.

В силу достаточности $S(\mathbf{X})$ выполнено равенство

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

Проверим, что $E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X}))$ равно тому же самому. Свойство измеримости уже есть, проверим интегральное. Для любого $A \in \mathcal{F}_{S(\mathbf{X})}$

$$E_{\theta, \theta'}(f_{\theta}(\mathbf{X}) \cdot I_A) = \frac{1}{2}(E_{\theta}(f_{\theta}(\mathbf{X})I_A) + E_{\theta'}(f_{\theta}(\mathbf{X}) \cdot I_A)) =$$

(интегральное свойство)

$$= \frac{1}{2}(E_{\theta}(E_{\theta}(f_{\theta}(\mathbf{X})|S(\mathbf{X})) \cdot I_A) + E_{\theta'}(E_{\theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X})) \cdot I_A))$$

$$= \frac{1}{2}(E_{\theta}(E_{\theta}(f_{\theta}(\mathbf{X})|S(\mathbf{X})) \cdot I_A) + E_{\theta'}(E_{\theta}(f_{\theta}(\mathbf{X})|S(\mathbf{X})) \cdot I_A)) = E_{\theta, \theta'}(E_{\theta}(f_{\theta}(\mathbf{X})|S(\mathbf{X})) \cdot I_A).$$

Далее, по формуле пересчета УМО

$$E_{\theta}(f_{\theta}(\mathbf{X})|S(\mathbf{X})) = \frac{E_{\theta, \theta'}(f_{\theta}^2(\mathbf{X})|S(\mathbf{X}))}{E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X}))}.$$

Но мы только что доказали, что $E_{\theta}(f_{\theta}(\mathbf{X})|S(\mathbf{X})) = E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X}))$, таким образом,

$$E_{\theta, \theta'}(f_{\theta}^2(\mathbf{X})|S(\mathbf{X})) = (E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X})))^2.$$

Вычитая из левой части правую и беря

$$E_{\theta, \theta'} f_{\theta}^2(\mathbf{X}) - E_{\theta, \theta'}(E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X})))^2 = E_{\theta, \theta'}(f_{\theta}(\mathbf{X}) - E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X})))^2 = 0.$$

Последнее означает, что $f_{\theta}(\mathbf{X}) = E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X}))$, т.е. она $S(\mathbf{X})$ -измерима.

\Rightarrow Из факторизации следует, что $f_{\theta}(\mathbf{X})$ и $f_{\theta'}(\mathbf{X})$ есть функции от $S(\mathbf{X})$. По формуле для пересчета УМО получаем, что для любой статистики $T(\mathbf{X})$ выполнено

$$E_{\theta}(T(\mathbf{X})|S(\mathbf{X})) = \frac{E_{\theta, \theta'}(T(\mathbf{X})f_{\theta}(\mathbf{X})|S(\mathbf{X}))}{E_{\theta, \theta'}(f_{\theta}(\mathbf{X})|S(\mathbf{X}))} = E_{\theta, \theta'}(T(\mathbf{X})|S(\mathbf{X}))$$

в силу $S(\mathbf{X})$ -измеримости $f_{\theta}(\mathbf{X})$. Совершенно аналогично, получаем такое же равенство для θ' :

$$E_{\theta'}(T(\mathbf{X})|S(\mathbf{X})) = \frac{E_{\theta, \theta'}(T(\mathbf{X})f'_{\theta}(\mathbf{X})|S(\mathbf{X}))}{E_{\theta, \theta'}(f'_{\theta}(\mathbf{X})|S(\mathbf{X}))} = E_{\theta, \theta'}(T(\mathbf{X})|S(\mathbf{X})).$$

Следовательно, для любого $\theta \in \Theta$ выполнено

$$E_{\theta}(T(\mathbf{X})|S(\mathbf{X})) = E_{\theta'}(T(\mathbf{X})|S(\mathbf{X})) \text{ почти наверное}$$

Взяв $T(\mathbf{x}) = I\{\mathbf{x} \in B\}$, получаем, что существует вариант условного распределения $P_{\theta}(\mathbf{X} \in B|S(\mathbf{X}))$, который не зависит от θ . По определению мы получаем, что $S(\mathbf{X})$ - достаточная статистика. ■

Примеры

Пример (1)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из распределения Бернулли $\text{Bin}(1, \theta)$, $\theta \in (0, 1)$. Найдите достаточную статистику.

Решение

Плотность имеет вид:

$$\begin{aligned} p_{\theta}(X_1, \dots, X_n) &= \prod_{i=1}^n p_{\theta}(X_i) = \prod_{i=1}^n \theta^{X_i} (1 - \theta)^{1 - X_i} \cdot I\{X_i \in \{0, 1\}\} = \\ &= \theta^{\sum_{i=1}^n X_i} (1 - \theta)^{n - \sum_{i=1}^n X_i} \cdot I\{\forall X_i \in \{0, 1\}\}. \end{aligned}$$

Значит, $\sum_{i=1}^n X_i$ - достаточная статистика.

Пример (2)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из пуассоновского распределения $\text{Pois}(\theta)$, $\theta > 0$.
Найдите достаточную статистику.

Решение

Плотность имеет вид:

$$p_{\theta}(X_1, \dots, X_n) = \prod_{i=1}^n p_{\theta}(X_i) = \prod_{i=1}^n \frac{\theta^{X_i}}{X_i!} e^{-\theta} = \left(\frac{1}{\prod_{i=1}^n X_i!} \right) \theta^{\sum_{i=1}^n X_i} \cdot e^{-n\theta}.$$

Значит, $\sum_{i=1}^n X_i$ - достаточная статистика.

Пример (3)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из нормального распределения $\mathcal{N}(a, \sigma^2)$, $\theta = (a, \sigma^2)$, $a \in \mathbb{R}$, $\sigma > 0$. Найдите достаточную статистику.

Решение

Плотность имеет вид:

$$p_{\theta}(X_1, \dots, X_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - a)^2}{2\sigma^2}\right) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{\sum_{i=1}^n X_i^2 - 2a \sum_{i=1}^n X_i + na^2}{2\sigma^2}\right).$$

Значит, $\left(\sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i\right)$ - достаточная статистика.

Пример (4)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из равномерного распределения $U(0, \theta)$, $\theta > 0$.
Найдите оптимальную оценку θ .

Решение

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

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

Достаточной статистикой в данном случае будет являться $X_{(n)}$, так как в критерии факторизации можно взять $h(\mathbf{X}) = I\{0 \leq X_{(1)} \leq X_{(n)}\}$, $\psi_{\theta}(S(\mathbf{X})) = \theta^{-n} I\{X_{(n)} \leq \theta\}$. Покажем теперь, что $X_{(n)}$ есть полная статистика. Для этого вспомним, что $P_{\theta}(X_{(n)} \leq x) = (\frac{x}{\theta})^n$ для $x \in [0, \theta]$, поэтому плотность максимума равна

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

Посчитаем математическое ожидание произвольной функции $g(X_{(n)})$:

$$E_{\theta} g(X_{(n)}) = \int_0^{\theta} g(x) \frac{nx^{n-1}}{\theta^n} dx.$$

Теперь предположим, что для всех $\theta > 0$ $E_{\theta} g(X_{(n)}) = 0$. Это равносильно тому, что

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

Продифференцируем по θ :

$$g(\theta) \theta^{n-1} = 0 \iff g(\theta) = 0 \text{ почти наверное для } \theta > 0.$$

Тогда $g(X_{(n)}) = 0$ P_{θ} - почти наверное для все $\theta > 0$ и $X_{(n)}$ есть полная достаточная статистика.

Осталось решить уравнение несмещённости: $E_{\theta} \varphi(X_{(n)}) = \theta$. Заметим, что это равносильно тому, что:

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

Продифференцируем по θ :

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

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

13.

Лекция 13

На этой лекции мы разберемся с тем, как можно проверять полноту достаточной статистики в ситуации, когда мы имеем дело с экспоненциальными семействами распределений, и начнем обсуждать доверительные интервалы (интервальный подход к оценке параметра).

Экспоненциальные семейства

Для достаточности есть общий критерий факторизации, который дает ответ в большинстве разумных случаев. Для полноты такого удобного критерия нет. Однако в некоторых хороших случаях можно вывести простые достаточные условия полноты - например, если семейство распределений является экспоненциальным.

Определение 13.1. Пусть $\{P_\theta, \theta \in \Theta\}$ - семейство распределений с k -мерным параметром: $\Theta \subseteq \mathbb{R}^k$, $k \in \mathbb{N}$. Далее, пусть оно является доминируемым относительно меры μ с плотностью $p_\theta(x)$. Если эта плотность представима в следующем виде:

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

где $h(\mathbf{x}) \geq 0$ и $U_i(\mathbf{x})$ - это борелевские функции, то семейство распределений $\{P_\theta, \theta \in \Theta\}$ называется экспоненциальным.

Теперь предположим, что семейство является экспоненциальным. Какая для него будет достаточная статистика? Понятно, что она будет равна $S(\mathbf{X}) = (U_1(\mathbf{X}), \dots, U_k(\mathbf{X}))$. Оказывается, что если функции $a_i\theta$ достаточно хороши, то $S(\mathbf{X})$ будет и полной.

Теорема 13.1. (об экспоненциальном семействе)

Пусть \mathbf{X} - наблюдение с неизвестным распределением, принадлежащим экспоненциальному семейству $\{P_\theta, \theta \in \Theta\}$ с плотностью вида (1). Если при пробегании θ всего Θ вектор $\mathbf{a}(\theta) = (a_1(\theta), \dots, a_k(\theta))$ зачерчивает k -мерный параллелепипед или шар (это значит, что внутри множества значений вектора $\mathbf{a}(\theta)$, при θ пробегающем Θ , содержится некоторый k -мерный параллелепипед или шар), то статистика $S(\mathbf{X}) = (U_1(\mathbf{X}), \dots, U_k(\mathbf{X}))$ будет полной достаточной статистикой для $\{P_\theta, \theta \in \Theta\}$.

Для доказательства потребуется несколько утверждений.

Плотность $S(\mathbf{X})$

Доказательство теоремы об экспоненциальном семействе использует комплексный анализ. Нам понадобится теорема единственности для аналитических функций. Но для начала разберемся с вопросом: чему равна плотность $S(\mathbf{X})$? Когда мы говорим о полноте, мы должны проверить, что математическое ожидание функции от $S(\mathbf{X})$ может быть нулем только если эта функция нулевая. Для того, чтобы выписать математическое ожидание в интегральном виде, нам необходимо понимать, а имеет ли $S(\mathbf{X})$ плотность. Оказывается, что плотность у $S(\mathbf{X})$ есть. Обозначим

$$\psi_\theta(S(\mathbf{X})) = \exp\left\{\sum_{i=1}^k U_i(\mathbf{X})a_i(\theta) + v(\theta)\right\} = \exp\{\langle \mathbf{a}(\theta), S(\mathbf{X}) \rangle + v(\theta)\}.$$

Лемма (1)

Статистика $S(\mathbf{X})$ будет иметь плотность $\psi_\theta(\mathbf{s})$ по мере ν на $(\mathbb{R}^k, \mathcal{B}(\mathbb{R}^k))$, заданной соотношением: для $B \in \mathcal{B}(\mathbb{R}^k)$

$$\nu(B) = \int_{S^{-1}(B)} h(\mathbf{x})\mu(d\mathbf{x}).$$

Доказательство.

Для существования плотности нужно показать, что распределение статистики $S(\mathbf{X})$ абсолютно непрерывно относительно меры ν . Пусть $G_\theta = P_\theta(S(\mathbf{X}) \in B)$ есть распределение $S(\mathbf{X})$. Если $\nu(B) = 0$, то $h(\mathbf{x}) = 0$ ν -п.н. на множестве интегрирования $S^{-1}(B)$. Тогда

$$G_\theta(B) = P_\theta(S(\mathbf{X}) \in B) = \int_{S^{-1}(B)} p_\theta(\mathbf{x})\mu(d\mathbf{x}) = \int_{S^{-1}(B)} h(\mathbf{x})\psi_\theta(S(\mathbf{x}))\mu(d\mathbf{x}) = 0.$$

Следовательно, распределение G_θ абсолютно непрерывно относительно ν и по теореме Радона-Никодима существует плотность $\frac{dG_\theta}{d\nu}$. Но в силу определения ν

$$G_\theta(B) = \int_{S^{-1}(B)} h(\mathbf{x})\psi_\theta(S(\mathbf{x}))\mu(d\mathbf{x}) = \int_B \psi_\theta(\mathbf{s})\nu(d\mathbf{s}).$$

Значит, $\frac{dG_\theta}{d\nu}(\mathbf{s}) = \psi_\theta(\mathbf{s})$.

■

Достаточные статистики и информация Фишера

Замечание

Лемма (1) выполнена не только для экспоненциального семейства, а всегда, когда имеет место факторизация.

Теперь мы можем доказать и последнее неразобранное свойство информации Фишера.

Следствие

В условиях регулярности $I_S(\theta) = I_X(\theta)$ для любого $\theta \in \Theta$ тогда и только тогда, когда статистика $S(\mathbf{X})$ является достаточной для семейства распределений $\{P_\theta, \theta \in \Theta\}$.

Доказательство.

\Rightarrow Пусть $S(\mathbf{X})$ - достаточная статистика. Тогда по критерию факторизации выполнено

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

где $\psi_\theta(\mathbf{s})$ - плотность $S(\mathbf{X})$. Стало быть, вклады \mathbf{X} и $S(\mathbf{X})$ равны:

$$\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) = \frac{\partial}{\partial \theta} \ln \psi_\theta(S(\mathbf{X})),$$

откуда следует и равенство информации Фишера, так как информация - второй момент для вклада.

\Leftarrow В обратную сторону почти все следует из доказательства, приведенного на лекции 5. Из равенства информации Фишера мы получаем, что вклады \mathbf{X} и $S(\mathbf{X})$ должны быть пропорциональны

$$\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{X}) = a(\theta) \cdot \frac{\partial}{\partial \theta} \ln g_\theta(S(\mathbf{X})),$$

где $g_\theta(\mathbf{s})$ - плотность $S(\mathbf{X})$. Находя из равенства $p_\theta(\mathbf{x})$, мы получаем факторизацию плотности:

$$\ln p_\theta(\mathbf{X}) = \ln \psi_\theta(S(\mathbf{X})) + h(\mathbf{X}).$$

Значит, $S(\mathbf{X})$ - достаточная статистика. ■

Преобразование Лапласа

Лемма (2)

Пусть G_1 и G_2 - две σ -конечные меры на $(\mathbb{R}^k, \mathcal{B}(\mathbb{R}^k))$. Если имеет место равенство конечных интегралов

$$\int_{\mathbb{R}^k} e^{\langle \mathbf{a}, \mathbf{u} \rangle} G_1(d\mathbf{u}) = \int_{\mathbb{R}^k} e^{\langle \mathbf{a}, \mathbf{u} \rangle} G_2(d\mathbf{u}) < \infty$$

для всех $\mathbf{a} \in \mathbf{I}$, где \mathbf{I} - некоторый параллелепипед в \mathbb{R}^k , то $G_1 = G_2$.

Доказательство.

Докажем только для $k = 1$ (при $k > 1$ все аналогично, но требует знаний многомерного комплексного анализа).

Сначала считаем, что $\mathbf{I} = \{|x| \leq \alpha\}$ для всех $\alpha > 0$. Тогда функции

$$h_i(a) = \int_{\mathbb{R}} e^{au} G_i(du), \quad i = 1, 2,$$

будут аналитическими в интервале $|a| < \alpha$. Рассмотрим теперь h_i как функции комплексного переменного:

$$h_i(z) = \int_{\mathbb{R}} e^{zu} G_i(du), \quad i = 1, 2, \quad z \in \mathbb{C}.$$

Тогда $h_i(z)$ будут аналитическими в полосе $|\operatorname{Re}(z)| < \alpha$ и совпадающими на интервале $\{|\operatorname{Re}(z)| < \alpha, \operatorname{Im}(z) = 0\}$. По теореме единственности $h_1(z) = h_2(z)$ во всей полосе и, в частности, на мнимой прямой $\operatorname{Re}(z) = 0$. Таким образом, для любого $t \in \mathbb{R}$ выполнено:

$$\int_{\mathbb{R}} e^{itu} G_1(du) = \int_{\mathbb{R}} e^{itu} G_2(du).$$

Видно, что $\int_{\mathbb{R}} G_1(du) = \int_{\mathbb{R}} G_2(du) < +\infty$, поэтому можно считать, что G_1 и G_2 - это вероятностные меры. Но тогда мы получаем, что совпадают их характеристические функции. Значит, $G_1 = G_2$.

Если же $\mathbf{I} = \{|x - \alpha_0| \leq \alpha\}$, то надо перейти к мерам

$$G_i^*(du) = e^{\alpha_0 u} G_i(du), \quad i = 1, 2.$$

■

Доказательство теоремы об экспоненциальном семействе

Доказательство.

Достаточность статистики $S(\mathbf{X})$ следует из критерия факторизации Неймана-Фишера. Докажем ее полноту. Пусть $\varphi(\mathbf{s})$ - такая борелевская функция из \mathbb{R}^k в \mathbb{R} , что $E_\theta \varphi(S(\mathbf{X})) = 0$ для всех $\theta \in \Theta$. Заметим, что по лемме (1) это эквивалентно тому, что для любого $\theta \in \Theta$

$$E_\theta \varphi(S(\mathbf{X})) = \int_{\mathbb{R}^k} \varphi(\mathbf{s}) \psi_\theta(\mathbf{s}) \nu(d\mathbf{s}) = 0.$$

Введем две функции: $\varphi^+(\mathbf{s}) = \max\{\varphi(\mathbf{s}), 0\}$ и $\varphi^- = \max\{-\varphi(\mathbf{s}), 0\}$. Тогда для любого $\theta \in \Theta$

$$\int_{\mathbb{R}^k} \varphi^+(\mathbf{s}) \psi_\theta(\mathbf{s}) \nu(d\mathbf{s}) = \int_{\mathbb{R}^k} \varphi^-(\mathbf{s}) \psi_\theta(\mathbf{s}) \nu(d\mathbf{s}).$$

Вспомнив, что $\psi_\theta(\mathbf{s}) = e^{\langle \mathbf{a}(\theta), \mathbf{s} \rangle + v(\theta)}$, мы получаем, что для любого $\theta \in \Theta$

$$\int_{\mathbb{R}^k} e^{\langle \mathbf{a}(\theta), \mathbf{s} \rangle} G_+(d\mathbf{s}) = \int_{\mathbb{R}^k} e^{\langle \mathbf{a}(\theta), \mathbf{s} \rangle} G_-(d\mathbf{s}),$$

где $G_\pm(d\mathbf{s}) = \varphi^\pm(\mathbf{s}) \nu(d\mathbf{s})$ - новые меры.

Из условия теоремы следует, что для всех \mathbf{b} из некоторого параллелепипеда в \mathbb{R}^k выполняется равенство

$$\int_{\mathbb{R}^k} e^{\langle \mathbf{b}, \mathbf{s} \rangle} G_+(d\mathbf{s}) = \int_{\mathbb{R}^k} e^{\langle \mathbf{b}, \mathbf{s} \rangle} G_-(d\mathbf{s}).$$

По лемме (2) это означает, что $G_+ = G_-$, то есть $\varphi(\mathbf{s}) = 0$ ν -п.н. По определению получаем, что $S(\mathbf{X})$ - это полная статистика для $\{P_\theta, \theta \in \Theta\}$. ■

Пример нахождения оптимальной оценки

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из нормального распределения $\mathcal{N}(\mu, \sigma^2)$, $\theta = (\mu, \sigma^2)$, $\theta \in \mathbb{R} \times \mathbb{R}_+$. Найдите оптимальную оценку θ .

Решение

План решения:

Необходимо найти полную достаточную статистику. Достаточность мы проверим по критерию факторизации, а полноту по теореме об экспоненциальном семействе. Останется решить уравнение несмещенности.

Начнем с критерия факторизации. Для этого выпишем плотность нашей выборки:

$$\begin{aligned} p_{\theta} &= \prod_{k=1}^n p_{\theta}(X_k) = \prod_{k=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(X_k - \mu)^2}{2\sigma^2}\right\} = \frac{1}{(2\pi)^{\frac{n}{2}} \sigma^n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{k=1}^n (X_k - \mu)^2\right\} = \\ &= \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{k=1}^n X_k^2 + \frac{\mu}{\sigma^2} \sum_{k=1}^n X_k - \frac{n\mu^2}{2\sigma^2} - \frac{n}{2} \ln \sigma^2\right\}. \end{aligned}$$

Тогда это будет экспоненциальное семейство распределений, где

$$\mathbf{a}(\theta) = \left(\frac{\mu}{2\sigma^2}, \frac{-1}{2\sigma^2}\right), \quad S(\mathbf{X}) = \left(\sum_{k=1}^n X_k, \sum_{k=1}^n X_k^2\right),$$

$$v(\theta) = \frac{-n\mu^2}{2\sigma^2} - n \ln \sigma, \quad h(\mathbf{X}) = (2\pi)^{-\frac{n}{2}}.$$

Заметим, что $\mathbf{a}(\theta)$ пробегает $\mathbb{R} \times \mathbb{R}_+$. Тогда по теореме $S(\mathbf{X})$ есть полная достаточная статистика.

Теперь надо решить уравнения несмещенности, причем их можно решать по-отдельности для μ и σ^2 .

Сразу скажем, что $E_{\theta} \bar{\mathbf{X}} = \mu$, 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 второй момент выборочного среднего, пользуясь тем, что $\bar{\mathbf{X}} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$:

$$E_{\theta}(\bar{\mathbf{X}})^2 = \mu^2 + \frac{\sigma^2}{n}.$$

Следовательно, оптимальной оценкой σ^2 будет

$$\frac{n}{n-1} \cdot S^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{\mathbf{X}})^2.$$

Подведем итог и сформулируем общую схему поиска оптимальных оценок.

Пусть требуется отыскать оптимальную оценку $\tau(\theta)$. Тогда

- 1) Находим достаточную статистику $S(\mathbf{X})$ с помощью критерия факторизации.
- 2) Проверяем ее на полноту с помощью теоремы об экспоненциальном семействе.
- 3) Вычисляем распределение $S(\mathbf{X})$.
- 4) Решаем уравнение несмещенности:

$$E_{\theta}\varphi(S(\mathbf{X})) = \tau(\theta).$$

Ответ: $\varphi(S(\mathbf{X}))$ - это оптимальная оценка $\tau(\theta)$.

Доверительное оценивание

Пре предыдущие темы нашего курса относились к так называемому "точечному" оцениванию параметров. В его рамках параметру θ сопоставляется оценка $\hat{\theta}(\mathbf{X})$, то есть точка, в некотором смысле близка к θ . Но насколько близко на самом деле находится $\hat{\theta}(\mathbf{X})$ к θ ? Например, нам бы хотелось указать такое ε , что с большой вероятностью выполнено $|\hat{\theta}(\mathbf{X}) - \theta| < \varepsilon$.

Более общая постановка вопроса: указать интервал, который покрывает θ с вероятностью не менее γ (обычно берут $\gamma = 0.95$ или 0.99). Величина γ называется уровнем доверия (доверительной вероятностью), она отражает "степень готовности мириться с ошибкой".

Важный момент: для выбранного γ следует выбрать наименее короткий интервал (в идеале - его длина должна стремиться к нулю с ростом числа наблюдений).

Доверительные интервалы

Пусть \mathbf{X} - наблюдение с неизвестным распределением $P \in \{P_{\theta}, \theta \in \Theta\}$, $\Theta \subset \mathbb{R}$.

Определение 13.2. Пара статистик $(T_1(\mathbf{X}), T_2(\mathbf{X}))$ называется доверительным интервалом уровня доверия γ для параметра θ , если для любого $\theta \in \Theta$ выполняется неравенство:

$$P_{\theta}(T_1(\mathbf{X}) < \theta < T_2(\mathbf{X})) \geq \gamma.$$

Если данная вероятность равна γ для всех $\theta \in \Theta$, то доверительный интервал называется точным.

Доверительные области

Если же параметр θ является многомерным, $\Theta \subset \mathbb{R}^k$, $k > 1$, то можно строить доверительные интервалы для его компонент θ_i или скалярных функций $\tau(\theta) \in \mathbb{R}$. В общем случае вводится понятие доверительной области.

Определение 13.3. Подмножество $S(\mathbf{X}) \subset \Theta$ называется доверительной областью уровня доверия γ для параметра θ , если для любого $\theta \in \Theta$ выполнено

$$P_{\theta}(\theta \in S(\mathbf{X})) \geq \gamma.$$

Вопрос: как строить доверительные интервалы и области?

Метод центральной статистики

Предположим, что нашлась такая одномерная функция $G(\mathbf{X}, \theta)$, что ее распределение не зависит от θ , а сама она зависит от θ . Подобная функция называется **центральной статистикой**.

Пусть распределение $G(\mathbf{X}, \theta)$ нам известно, и для заданных $\gamma_1, \gamma_2 \in (0, 1)$ с условием $\gamma_2 - \gamma_1 = \gamma$ мы смогли вычислить величины g_i , $i = 1, 2$, - γ_i -квантили распределения $G(\mathbf{X}, \theta)$. Тогда для любого $\theta \in \Theta$ мы получаем, что

$$P_{\theta}(g_1 \leq G(\mathbf{X}, \theta) \leq g_2) \geq \gamma_2 - \gamma_1 = \gamma.$$

Утверждение

Множество $S(\mathbf{X}) = \{\theta \in \Theta : g_1 \leq G(\mathbf{X}, \theta) \leq g_2\}$ является доверительной областью уровня доверия γ для параметра θ .

Доказательство.

$$P_{\theta}(\theta \in S(\mathbf{X})) = P_{\theta}(g_1 \leq G(\mathbf{X}, \theta) \leq g_2) \geq \gamma.$$

■

Идеи по вычислению

Указанный подход имеет сразу несколько проблем.

- Неравенство $g_1 \leq G(\mathbf{X}, \theta) \leq g_2$ надо решить относительно θ . Если $G(\mathbf{X}, \theta)$ имеет абсолютно непрерывное распределение а сама монотонно и непрерывно зависит от θ , то $\{g_1 < G(\mathbf{X}, \theta) < g_2\}$ станет интервалом вида $\{T_1(\mathbf{X}) < \theta < T_2(\mathbf{X})\}$. Но если зависимость от θ "плохая", то может быть достаточно трудно решить данное неравенство относительно θ .

- Поиск центральной статистики - отдельная и сложная задача. Она легко решается, если проста природа параметра θ . Например, это параметр сдвига или масштаба.
- Даже для параметра сдвига или масштаба необходимо вычислить распределение центральной статистики $G(\mathbf{X}, \theta)$ и добиваться малой длины получающегося интервала.

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из равномерного распределения $U(0, \theta)$, $\theta > 0$. Постройте точный доверительный интервал уровня доверия γ с помощью статистики $X_{(n)}$.

Решение

Ясно, что $\frac{X_i}{\theta} \sim U(0, 1)$. Тогда

$$P_{\theta}\left(\frac{X_{(n)}}{\theta} \leq x\right) = x^n, \quad x \in [0, 1].$$

Значит,

$$P_{\theta}\left(g_1 < \frac{X_{(n)}}{\theta} < g_2\right) = g_2^n - g_1^n = \gamma.$$

Положим $g_2 = 1$, $g_1 = (1 - \gamma)^{\frac{1}{n}}$. Тогда с вероятностью в точности γ

$$X_{(n)} < \theta < X_{(n)}(1 - \gamma)^{-\frac{1}{n}}.$$

Длина интервала имеет порядок $O(\frac{1}{n})$.

Асимптотические доверительные интервалы

Вопрос: что делать, если центральную статистику не удастся быстро найти или ее распределение сложно вычислить?

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка растущего размера с неизвестным распределением $P \in \{P_{\theta}, \theta \in \Theta\}$, $\Theta \subset \mathbb{R}$.

Определение 13.4. Последовательность пар статистик $(T_n^1(X_1, \dots, X_n), T_n^2(X_1, \dots, X_n))$ называется асимптотическим доверительным интервалом уровня доверия γ для параметра θ , если для любого $\theta \in \Theta$

$$\lim_{n \rightarrow +\infty} P_{\theta}(T_n^1(X_1, \dots, X_n) < \theta < T_n^2(X_1, \dots, X_n)) \geq \gamma.$$

Если для любого $\theta \in \Theta$

$$\exists \lim_{n \rightarrow +\infty} P_{\theta}(T_n^1(X_1, \dots, X_n) < \theta < T_n^2(X_1, \dots, X_n)) = \gamma,$$

то асимптотический доверительный интервал называется точным.

Метод построения

Идея построения использует асимптотически нормальные оценки. Свойство асимптотической нормальности оказывается здесь принципиально важным, так как оно позволяет нам строить весьма хорошие асимптотические доверительные интервалы, которые часто используются на практике.

Пусть $\hat{\theta}_n(\mathbf{X}) = \hat{\theta}_n(X_1, \dots, X_n)$ - это асимптотически нормальная оценка θ с асимптотической дисперсией $\sigma^2(\theta)$, то есть для любого $\theta \in \Theta$

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

Тогда

$$\sqrt{n} \cdot \frac{\hat{\theta}_n(\mathbf{X}) - \theta}{\sigma(\theta)} \xrightarrow{d_{\theta}} \mathcal{N}(0, 1).$$

Надо избавиться от θ в знаменателе, заменив $\sigma(\theta)$ на ее состоятельную оценку $\hat{\delta}_n(\mathbf{X})$.

Тогда

$$\sqrt{n} \cdot \frac{\hat{\theta}_n(\mathbf{X}) - \theta}{\hat{\delta}_n(\mathbf{X})} = \sqrt{n} \cdot \frac{\hat{\theta}_n(\mathbf{X}) - \theta}{\sigma(\theta)} \cdot \frac{\sigma(\theta)}{\hat{\delta}_n(\mathbf{X})}.$$

Первый множитель сходится по распределению к $\mathcal{N}(0, 1)$, а второй - к 1 по вероятности. По лемме Слущкого получаем, что

$$\sqrt{n} \cdot \frac{\hat{\theta}_n(\mathbf{X}) - \theta}{\hat{\delta}_n(\mathbf{X})} \xrightarrow{d_{\theta}} \mathcal{N}(0, 1).$$

Теперь мы можем построить асимптотические доверительные интервалы. Обозначим через u_{α} - α -квантиль $\mathcal{N}(0, 1)$. В силу симметрии будет выполнено

$$u_{\frac{1+\gamma}{2}} = -u_{\frac{1-\gamma}{2}}.$$

Тогда для любого $\theta \in \Theta$

$$P_{\theta} \left(\sqrt{n} \cdot \left| \frac{\hat{\theta}_n(\mathbf{X}) - \theta}{\hat{\delta}_n(\mathbf{X})} \right| < u_{\frac{1+\gamma}{2}} \right) \rightarrow \gamma \quad \text{при } n \rightarrow +\infty.$$

Следовательно, получаем следующий асимптотический доверительный интервал:

$$\left(\hat{\theta}_n(\mathbf{X}) - u_{\frac{1+\gamma}{2}} \cdot \frac{\hat{\delta}_n(\mathbf{X})}{\sqrt{n}}, \hat{\theta}_n(\mathbf{X}) + u_{\frac{1+\gamma}{2}} \cdot \frac{\hat{\delta}_n(\mathbf{X})}{\sqrt{n}} \right).$$

Его длина имеет порядок $O(\frac{1}{\sqrt{n}})$.

Замечание (1)

Длина получаемого интервала равна

$$2u_{\frac{1+\gamma}{2}} \cdot \frac{\hat{\delta}_n(\mathbf{X})}{\sqrt{n}} \sim 2u_{\frac{1+\gamma}{2}} \cdot \frac{\sigma(\theta)}{\sqrt{n}},$$

то есть выбор оценки с меньшей асимптотической дисперсией уменьшает длину асимптотического доверительного интервала.

Замечание (2)

Как находить $\hat{\delta}_n(\mathbf{X})$? Заметим, что $\hat{\theta}_n(\mathbf{X}) \xrightarrow{P_\theta} \theta$. Если $\sigma(\theta)$ непрерывна по θ , то подойдет $\hat{\delta}_n(\mathbf{X}) = \sigma(\hat{\theta}_n(\mathbf{X}))$, которая является состоятельной оценкой $\sigma(\theta)$ в силу теоремы о наследовании сходимости.

Однако, бывает, что $\sigma(\theta)$ не удастся явно найти как функцию от параметра. В этом случае может подойти состоятельная оценка "общего вида". Например, если мы знаем, что $\sigma^2(\theta)$ - это дисперсия X_1 , то ее состоятельно оценивает выборочная дисперсия S^2 .

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из схемы Бернулли $\text{Bin}(1, \theta)$, $\theta \in (0, 1)$. Постройте асимптотический доверительный интервал уровня доверия γ .

Решение

В силу того, что

$$E_\theta X_1 = \theta, D_\theta X_1 = \theta(1 - \theta),$$

оценка $\bar{\mathbf{X}}$ будет асимптотически нормальной оценкой θ с асимптотической дисперсией $\sigma^2(\theta) = \theta(1 - \theta)$. Согласно ЦПТ

$$\sqrt{n} \cdot \frac{\bar{\mathbf{X}} - \theta}{\sqrt{\theta(1 - \theta)}} \xrightarrow{d_\theta} \mathcal{N}(0, 1).$$

Статистика $S^2 = \bar{\mathbf{X}}^2 - (\bar{\mathbf{X}})^2$ является состоятельной оценкой $\theta(1 - \theta)$, поэтому используя описанную схему, получаем следующий доверительный интервал:

$$\left(\bar{\mathbf{X}} - u_{\frac{1+\gamma}{2}} \cdot \frac{S}{\sqrt{n}}, \bar{\mathbf{X}} + u_{\frac{1+\gamma}{2}} \cdot \frac{S}{\sqrt{n}} \right).$$

14.

Лекция 14

Постановка задачи линейной регрессии

Разберем классическую задачу линейной регрессии. Она формулируется следующим образом: допустим, что мы хотим узнать какую-то величину $\mathbf{l} \in \mathbb{R}^n$ и проводим ее измерение. Но реальные инструменты не могут иметь абсолютную точность, поэтому замеры будут содержать ошибку:

$$\mathbf{X} = \mathbf{l} + \boldsymbol{\varepsilon},$$

где \mathbf{X} - результат измерения, а $\boldsymbol{\varepsilon}$ - ошибка измерения. Обычно предполагается, что $\boldsymbol{\varepsilon}$ есть гауссовский вектор, но сначала мы будем рассматривать общий случай. На $\boldsymbol{\varepsilon}$ накладываются следующие ограничения:

$$E\boldsymbol{\varepsilon} = \mathbf{0}, \quad D\boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}_n,$$

где \mathbf{I}_n - единичная матрица размера $n \times n$. В среднем ошибка нулевая, а между собой ошибки измерений разных компонент некоррелированы и имеют одинаковую известную дисперсию σ^2 . В таком виде можно представить очень многие реальные задачи.

Задача

По наблюдению \mathbf{X} оценить неизвестные параметры \mathbf{l} и σ^2 .

В чем состоит линейность модели? В том, что мы кое-что знаем о векторе \mathbf{l} : он принадлежит нетривиальному известному линейному подпространству L размерности $k < n$. Но что означает, что мы знаем L ?

Будем считать, что нам известен его некоторый базис $z_1, \dots, z_k \in \mathbb{R}^n$. Составим из базисных векторов регрессионную матрицу: $\mathbf{Z} = (z_1, \dots, z_k) \in \text{Mat}(n \times k)$. Понятно, что в таком случае $\mathbf{Z} = k$. Тогда \mathbf{l} можно разложить в линейную комбинацию:

$$\mathbf{l} = \sum_{i=1}^k \theta_i z_i = \mathbf{Z}\boldsymbol{\theta}, \quad \text{где } \boldsymbol{\theta} = (\theta_1, \dots, \theta_k)^T.$$

Смысл вектора $\boldsymbol{\theta}$ очевиден: это неизвестные координаты вектора \mathbf{l} в базисе z_1, \dots, z_k . Тогда задача линейной регрессии переформулируется следующим образом:

Переформулированная задача

По наблюдению \mathbf{X} нужно оценить параметры $\boldsymbol{\theta}$ и σ^2 .

Пример

Пусть объект движется по прямой с постоянной скоростью. Далее, в некоторые моменты времени t_1, \dots, t_n мы измеряем положение объекта:

$$X_i = a + bt_i + \varepsilon_i, \quad i = 1, \dots, n,$$

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 . Легко видеть, что это задача линейной регрессии, где

$$\mathbf{l} = \begin{pmatrix} a + bt_1 \\ \vdots \\ a + bt_n \end{pmatrix} = \begin{pmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_n \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \mathbf{Z}\boldsymbol{\theta},$$

$\boldsymbol{\theta} = (a, b)^T$. Можно провести аналогичные рассуждения, если рассматривать объект, который движется по параболе. В таком случае будет добавляться член квадратичной зависимости от времени, ct_i^2 .

Метод наименьших квадратов

Как было сказано ранее, задача состоит в оценивании пары $(\boldsymbol{\theta}, \sigma^2)$. Для этого максимально удобным оказывается метод наименьших квадратов.

Определение 14.1. Оценкой $\boldsymbol{\theta}$ по методу наименьших квадратов называется

$$\hat{\boldsymbol{\theta}}(\mathbf{X}) = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^n} \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2.$$

Сама $\hat{\boldsymbol{\theta}}$ называется оценкой наименьших квадратов (о.н.к.).

Геометрический смысл этого определения крайне прост: $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть проекция \mathbf{X} на L . Но тогда $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть проекция \mathbf{X} на ортогональное дополнение L^\perp . Отсюда можно получить явную формулу для $\hat{\boldsymbol{\theta}}(\mathbf{X})$

Лемма (1)

Оценка наименьших квадратов имеет следующий вид:

$$\hat{\boldsymbol{\theta}}(\mathbf{X}) = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X}.$$

Замечание

Заметим, что $\mathbf{Z}^T \in \text{Mat}(k \times k)$ является обратимой матрицей, так как столбцы \mathbf{Z} являются базисными.

Доказательство.

Как было сказано выше, $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть проекция \mathbf{X} на L^\perp . Следовательно, для любого $\boldsymbol{\theta} \in \mathbb{R}^k$ вектор $\mathbf{Z}\boldsymbol{\theta}$ ортогонален $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$:

$$\langle \mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X}), \mathbf{Z}\boldsymbol{\theta} \rangle = \boldsymbol{\theta}^T \mathbf{Z}^T (\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})) = 0$$

Но второе равенство выполнено для всех $\boldsymbol{\theta} \in \mathbb{R}^k$ тогда и только тогда, когда $\mathbf{Z}^T (\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})) = 0$. Следовательно,

$$\mathbf{Z}^T \mathbf{X} = \mathbf{Z}^T \mathbf{Z} \hat{\boldsymbol{\theta}}(\mathbf{X}) \implies \hat{\boldsymbol{\theta}}(\mathbf{X}) = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X}.$$

■

Свойства оценки наименьших квадратов

Попытаемся понять, какие свойства могут быть у оценок методом наименьших квадратов. Начнем с самого простого: посчитаем среднее и дисперсию.

Лемма (2)

Для оценки наименьших квадратов $\hat{\boldsymbol{\theta}}(\mathbf{X})$ выполнено следующее:

$$\mathbb{E} \hat{\boldsymbol{\theta}}(\mathbf{X}) = \boldsymbol{\theta}, \quad D \hat{\boldsymbol{\theta}}(\mathbf{X}) = \sigma^2 (\mathbf{Z}^T \mathbf{Z})^{-1}.$$

Доказательство.

Для начала заметим, что матрица $\mathbf{Z}^T \mathbf{Z}$ симметрична. Далее, $\mathbf{X} = \mathbf{Z}\boldsymbol{\theta} + \boldsymbol{\varepsilon}$, поэтому $\mathbb{E} \mathbf{X} = \mathbf{Z}\boldsymbol{\theta}$ и $D \mathbf{X} = \sigma^2 \mathbf{I}_n$. Тогда

$$\mathbb{E} \hat{\boldsymbol{\theta}}(\mathbf{X}) = \mathbb{E} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbb{E} \mathbf{X} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Z} \boldsymbol{\theta} = \boldsymbol{\theta},$$

$$D \hat{\boldsymbol{\theta}}(\mathbf{X}) = D [(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X}] = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T D \mathbf{X} \cdot \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} = \sigma^{-1} (\mathbf{Z}^T \mathbf{Z})^{-1}.$$

■

Лемма (3)

Оценка методом наименьших квадратов $\hat{\boldsymbol{\theta}}(\mathbf{X})$ является оптимальной оценкой $\boldsymbol{\theta}$ в классе линейных несмещенных оценок, то есть оценок вида $\mathbf{B}\mathbf{X}$, где \mathbf{B} - неслучайная матрица.

Доказательство.

Пусть $\boldsymbol{\theta}^*(\mathbf{X}) = \mathbf{B}\mathbf{X}$ - другая несмещенная оценка $\boldsymbol{\theta}$. Тогда для любого $\boldsymbol{\theta} \in \mathbb{R}^k$

$$\boldsymbol{\theta} = \mathbb{E}_{\boldsymbol{\theta}} \boldsymbol{\theta}^*(\mathbf{X}) = \mathbf{B} \mathbb{E}_{\boldsymbol{\theta}} \mathbf{X} = \mathbf{B} \mathbf{Z} \boldsymbol{\theta}.$$

Но $\mathbf{B} \in \text{Mat}(k \times n)$, $\mathbf{Z} \in \text{Mat}(n \times k)$ и $\mathbf{B} \mathbf{Z} \in \text{Mat}(k \times k)$. Стало быть, $\mathbf{B} \mathbf{Z} = \mathbf{I}_k$. Далее заметим, что

$$D_{\boldsymbol{\theta}} \mathbf{B} \mathbf{X} = \mathbf{B} D_{\boldsymbol{\theta} \mathbf{X} \mathbf{B}^T = \sigma^2 \mathbf{B} \mathbf{B}^T}.$$

В итоге нужно доказать, что $\sigma^2 \mathbf{B} \mathbf{B}^T \geq \sigma^2 (\mathbf{Z}^T \mathbf{Z})^{-1}$. Для этого рассмотрим следующую дисперсию:

$$\begin{aligned} D_{\theta}[(\mathbf{B} - (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X})] &= (\mathbf{B} - (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T) \sigma^2 \mathbf{I}_n (\mathbf{B}^T - \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1}) = \\ &= \sigma^2 [\mathbf{B} \mathbf{B}^T - \mathbf{B} \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} - (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{B}^T + (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1}] = \\ &= \sigma^2 \mathbf{B} \mathbf{B}^T - \sigma^2 (\mathbf{Z}^T \mathbf{Z})^{-1} \geq 0. \end{aligned}$$

■

Оценка σ^2

Отлично, мы получили наилучшую линейную оценку для θ . Но еще есть неизвестный параметр σ^2 , для которого можно предложить следующую хорошую оценку.

Лемма (4)

Пусть $\hat{\theta}(\mathbf{X})$ - это оценка наименьших квадратов. Тогда

$$E_{\theta} \|\mathbf{X} - \mathbf{Z} \hat{\theta}(\mathbf{X})\|^2 = (n - k) \sigma^2.$$

Тем самым получаем, что несмещенной оценкой σ^2 является

$$\hat{\sigma}^2(\mathbf{X}) = \frac{1}{n - k} \|\mathbf{X} - \mathbf{Z} \hat{\theta}(\mathbf{X})\|^2.$$

Доказательство.

Для начала заметим, что $E_{\theta}(\mathbf{X} - \mathbf{Z} \hat{\theta}(\mathbf{X})) = 0$. Тогда

$$\begin{aligned} E_{\theta} \|\mathbf{X} - \mathbf{Z} \hat{\theta}(\mathbf{X})\|^2 &= \sum_{i=1}^n D_{\theta}(X_i - (\mathbf{Z} \hat{\theta}(\mathbf{X}))_i) = \text{tr } D_{\theta}(\mathbf{X} - \mathbf{Z} \hat{\theta}(\mathbf{X})) = \\ &= \text{tr } D_{\theta}(\mathbf{I}_n - \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T) \mathbf{X} = \sigma^2 \text{tr } (\mathbf{I}_n - \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T) (\mathbf{I}_n - \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T) = \\ &= \sigma^2 \text{tr}(\mathbf{I}_n - 2 \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T + \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T) = \\ &= \sigma^2 \text{tr}(\mathbf{I}_n - \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T) = \sigma^2 (n - \text{tr}(\mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T)) = \\ &= \sigma^2 (n - \text{tr}((\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Z})) = \sigma^2 (n - k). \end{aligned}$$

■

Гауссовская линейная модель

Гауссовская линейная модель - это линейная регрессионная модель

$$\mathbf{X} = \mathbf{l} + \boldsymbol{\varepsilon},$$

в которой

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n) \implies \mathbf{X} \sim \mathcal{N}(\mathbf{l}, \sigma^2 \mathbf{I}_n) = \mathcal{N}(\mathbf{Z}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_n).$$

Неизвестные параметры: $\boldsymbol{\theta} \in \mathbb{R}^k, \sigma > 0$.

Оказывается, подобное распределение $\boldsymbol{\varepsilon}$ позволяет решить сразу несколько задач:

- найти оптимальные оценки параметров;
- вычислить их распределения;
- построить точные доверительные интервалы и области.

Достаточные статистики

Пусть $\hat{\boldsymbol{\theta}}(\mathbf{X})$ - оценка наименьших квадратов.

Теорема 14.1. В гауссовской линейной регрессионной модели пара $(\hat{\boldsymbol{\theta}}(\mathbf{X}), \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2)$ является полной достаточной статистикой.

Доказательство.

Найдем плотность \mathbf{X} . Вспомним, что $\mathbf{X} \sim \mathcal{N}(\mathbf{l}, \sigma^2 \mathbf{I}_n)$, то есть компоненты \mathbf{X} независимы. Тогда

$$\begin{aligned} p(\mathbf{X}) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(X_i - l_i)^2\right\} = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left\{-\frac{1}{2\sigma^2}\|\mathbf{X} - \mathbf{l}\|^2\right\} = \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left\{-\frac{1}{2\sigma^2}\|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2\right\} = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left\{-\frac{1}{2\sigma^2}(\|\mathbf{X}\|^2 + \|\mathbf{Z}\boldsymbol{\theta}\|^2 - 2\boldsymbol{\theta}^T \mathbf{Z}^T \mathbf{X})\right\} = \\ &= \exp\left\{-\frac{\|\mathbf{X}\|^2}{2\sigma^2} + \sum_{i=1}^k \frac{\theta_i}{\sigma^2} (\mathbf{Z}^T \mathbf{X})_i - \frac{\|\mathbf{Z}\boldsymbol{\theta}\|^2}{2\sigma^2} - \frac{n}{2} \cdot \ln(2\pi\sigma^2)\right\}. \end{aligned}$$

При перебегании $\boldsymbol{\theta} \in \mathbb{R}^k, \sigma > 0$ вектор $(\frac{\theta_1}{\sigma^2}, \dots, \frac{\theta_k}{\sigma^2}, -\frac{1}{2\sigma^2})$ зачертит подпространство \mathbb{R}^{k+1} . По теореме об экспоненциальном семействе вектор $(\mathbf{Z}^T \mathbf{X}, \|\mathbf{X}\|^2)$ является полной достаточной статистикой.

Построим биекцию между $(\mathbf{Z}^T \mathbf{X}, \|\mathbf{X}\|^2)$ 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

$$\mathbf{Z}^T \mathbf{X} = (\mathbf{Z}^T \mathbf{Z}) \hat{\boldsymbol{\theta}}(\mathbf{X}).$$

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 $(\mathbf{Z}^T \mathbf{X}, \|\mathbf{X}\|^2)$ (так как \mathbf{Z} тоже известна).

В обратную сторону все аналогично:

$$\hat{\boldsymbol{\theta}}(\mathbf{X}) = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X},$$

$$\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 = \|\mathbf{X}\|^2 - \|\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2.$$

Тем самым получаем полную биективную измеримую связь между парами. Следовательно, $(\hat{\boldsymbol{\theta}}(\mathbf{X}), \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2)$ тоже есть полная достаточная статистика. ■

Следствие

В линейной гауссовской модели:

- 1) $\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть оптимальная оценка $\boldsymbol{\theta}$,
- 2) $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть оптимальная оценка \mathbf{l} ,
- 3) $\frac{1}{n-k} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2$ есть оптимальная оценка σ^2 .

Доказательство.

Ранее было показано, что эти величины являются несмещенными оценками для соответствующих параметров. Но они все являются функциями от полной и достаточной статистики. Следовательно, они оптимальны. ■

Доверительные интервалы и области

Теперь мы бы хотели построить доверительные интервалы и области для параметров линейной регрессии. Для этого нам понадобится теорема об ортогональных разложениях, но перед этим вспомним, что такое распределение хи-квадрат.

Определение 14.2. Будем говорить, что случайная величина ξ имеет распределение хи-квадрат с n степенями свободы, если ее плотность равна

$$p_{\xi}(x) = \frac{x^{\frac{n}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} \mathbf{I}\{x \geq 0\},$$

то есть $\xi \sim \Gamma(\frac{1}{2}, \frac{n}{2})$. Обозначение: $\xi \sim \chi_n^2$.

Упражнение

Если ξ_1, \dots, ξ_n есть независимые одинаково распределенные $\mathcal{N}(0, 1)$ случайные величины, то $\xi_1^2 + \dots + \xi_n^2 \sim \chi_n^2$.

Теорема 14.2. (об ортогональных разложениях) Пусть $\mathbf{X} \sim \mathcal{N}(\mathbf{l}, \sigma^2 \mathbf{I}_n)$, 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 $\mathbf{Y}_1, \dots, \mathbf{Y}_r$ независимы в совокупности и

$$\frac{1}{\sigma^2} \|\mathbf{Y}_i - \mathbf{EY}_i\|^2 \sim \chi_{\dim L_i}^2.$$

Доказательство.

Идея состоит в том, чтобы свести эту теорему к факту, сформулированному выше. Возьмем в \mathbb{R}^n такой ортонормированный базис $\mathbf{f}_1, \dots, \mathbf{f}_n$, что

$\mathbf{f}_1, \dots, \mathbf{f}_{k_1}$ - базис L_1 ,

$\mathbf{f}_{k_1+1}, \dots, \mathbf{f}_{k_1+k_2}$ - базис L_2 ,

\vdots

$\mathbf{f}_{k_1+\dots+k_{r-1}+1}, \dots, \mathbf{f}_n$ - базис L_r .

В таком случае размерность L_i равна k_i для всех $i = 1, \dots, r$. Далее, введем коэффициенты $W_i = \langle \mathbf{X}, \mathbf{f}_i \rangle$ для всех $i = 1, \dots, n$. Соберем эти коэффициенты в вектор и представим его в виде линейного преобразования \mathbf{X} :

$$\mathbf{W} = \begin{pmatrix} W_1 \\ \vdots \\ W_n \end{pmatrix} = \mathbf{C}\mathbf{X}, \text{ где } \mathbf{C} = \begin{pmatrix} \mathbf{f}_1^T \\ \vdots \\ \mathbf{f}_n^T \end{pmatrix}.$$

Заметим, что матрица \mathbf{C} ортогональна, так как ее строки задают ортонормированный базис \mathbb{R}^n . Вектор \mathbf{X} является гауссовским, значит, и \mathbf{W} тоже будет гауссовским, так как он получается линейным преобразованием \mathbf{X} .

Найдем распределение \mathbf{W} :

$$\mathbf{E}\mathbf{W} = \mathbf{E}\mathbf{C}\mathbf{X} = \mathbf{C}\mathbf{l};$$

$$\mathbf{D}\mathbf{W} = \mathbf{C} \cdot \mathbf{D}\mathbf{X} \cdot \mathbf{C}^T = \sigma^2 \mathbf{C}\mathbf{C}^T = \sigma^2 \mathbf{I}_n.$$

Следовательно, $\mathbf{W} \sim \mathcal{N}(\mathbf{Cl}, \sigma^2 \mathbf{I}_n)$ и, значит, компоненты \mathbf{W} независимы в совокупности.

Далее, случайные векторы \mathbf{Y}_j , которые вычисляются 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 найдем искомые распределения их длин. 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, в силу ортонормированности базиса \mathbf{f}_j получаем, что

$$\begin{aligned} \frac{1}{\sigma^2} \|\mathbf{Y}_j - \mathbb{E}\mathbf{Y}_j\|^2 &= \left\| \sum_{i=1}^{k_j} \frac{W_{k_1+\dots+k_{j-1}+i} - \mathbb{E}W_{k_1+\dots+k_{j-1}+i}}{\sigma} \mathbf{f}_{k_1+\dots+k_{j-1}+i} \right\|^2 = \\ &= \sum_{i=1}^{k_j} \left(\frac{W_{k_1+\dots+k_{j-1}+i} - \mathbb{E}W_{k_1+\dots+k_{j-1}+i}}{\sigma} \right)^2 \sim \chi_{k_j}^2. \end{aligned}$$

■

Распределение статистик

Сформулируем очень важное следствие:

Следствие

В гауссовской линейной модели оценка $\hat{\boldsymbol{\theta}}(\mathbf{X})$ независима с $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ и

$$\frac{1}{\sigma^2} \|\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X}) - \mathbf{Z}\boldsymbol{\theta}\|^2 \sim \chi_k^2,$$

$$\frac{1}{\sigma^2} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 \sim \chi_{n-k}^2.$$

Доказательство.

Мы знаем, что $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть проекция \mathbf{X} на L , а $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ проецирует \mathbf{X} на L^\perp . Согласно теореме об ортогональных разложениях $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ и $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$ независимы. В силу того, что $\mathbb{E}(\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})) = \mathbf{0}$, $\mathbb{E}\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X}) = \mathbf{Z}\boldsymbol{\theta}$, мы автоматически получаем и утверждения о распределениях.

Осталось заметить, что $\hat{\boldsymbol{\theta}}(\mathbf{X})$ есть функция от $\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$:

$\hat{\boldsymbol{\theta}}(\mathbf{X}) = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T (\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X}))$. Значит, $\hat{\boldsymbol{\theta}}(\mathbf{X})$ независима с $\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})$.

■

1. Доверительный интервал для σ^2

Воспользуемся методом центральной статистики. Заметим, что в данном случае можно использовать $\sigma^{-2} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2$ в качестве центральной статистики, так как

ранее было доказано, что эта величина имеет распределение хи-квадрат с $n - k$ степенями свободы.

Пусть c_α - это α -квантиль распределения χ^2_{n-k} . Тогда

$$P\left(c_{\frac{1-\gamma}{2}} < \frac{1}{\sigma^2} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 < c_{\frac{1+\gamma}{2}}\right) = \gamma.$$

Следовательно, точный доверительный интервал для σ^2 уровня доверия γ имеет вид:

$$\left(\frac{1}{c_{\frac{1+\gamma}{2}}} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2, \frac{1}{c_{\frac{1-\gamma}{2}}} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2\right).$$

2. Доверительный интервал для θ_i

Определение 14.3. Пусть с.в. ξ и η независимы, $\xi \sim \mathcal{N}(0, 1)$, $\eta \sim \chi^2_k$. Тогда случайная величина

$$\zeta = \frac{\xi}{\sqrt{\frac{\eta}{k}}}$$

имеет распределение Стьюдента с k степенями свободы (Т-распределение). Обозначение: $\zeta \sim T_k$.

Построим доверительный интервал для θ_i . Для этого заметим, что $\hat{\boldsymbol{\theta}}(\mathbf{X}) \sim \mathcal{N}(\boldsymbol{\theta}, \sigma^2 \mathbf{A})$, где $\mathbf{A} = (\mathbf{Z}^T \mathbf{Z})$. Следовательно,

$$\frac{\hat{\theta}_i - \theta_i}{\sqrt{\sigma^2 a_{ii}}} \sim \mathcal{N}(0, 1).$$

Далее, $\hat{\theta}_i$ не зависит от $\|\mathbf{X} - \mathbf{Z}\|$, а $\sigma^{-2} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 \sim \chi^2_{n-k}$. Значит, величина

$$\sqrt{\frac{n-k}{a_{ii}}} \frac{\hat{\theta}_i - \theta_i}{\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|} \sim T_{n-k}$$

имеет распределение Стьюдента T_{n-k} .

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:

$$\theta_i \in \left(\hat{\theta}_i - t_{\frac{1+\gamma}{2}} \sqrt{\frac{a_{ii}}{n-k}} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|, \hat{\theta}_i + t_{\frac{1+\gamma}{2}} \sqrt{\frac{a_{ii}}{n-k}} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|\right).$$

3. Доверительный интервал для θ

Осталось построить доверительную область для θ . Здесь возникает третье стандартное распределение, связанное с гауссовским - распределение Фишера.

Определение 14.4. Пусть ξ, η - независимые случайные величины, причем $\xi \sim \chi_k^2, \eta \sim \chi_n^2$. Тогда случайная величина $\zeta = \frac{\xi/k}{\eta/n}$ имеет распределение Фишера (распределение Снедекора, F-распределение) с (k, n) степенями свободы. Обозначение $\zeta \sim F_{k,n}$.

Согласно следствию из теоремы об ортогональных разложениях мы знаем, что

$$\frac{1}{\sigma^2} \|\mathbf{Z}\hat{\theta}(\mathbf{X}) - \mathbf{Z}\theta\|^2 \sim \chi_k^2,$$

$$\frac{1}{\sigma^2} \|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2 \sim \chi_{n-k}^2,$$

и эти случайные величины независимы.

Значит,

$$\frac{\|\mathbf{Z}\hat{\theta}(\mathbf{X}) - \mathbf{Z}\theta\|^2 n - k}{\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2 k} \sim F_{k, n-k}.$$

Далее, пусть f_α есть α -квантиль распределения $F_{k, n-k}$. Так как оно неотрицательно, то доверительной областью для всего вектора θ будет

$$S(\mathbf{X}) = \left\{ \theta \in \mathbb{R}^k : \frac{\|\mathbf{Z}\hat{\theta}(\mathbf{X}) - \mathbf{Z}\theta\|^2 n - k}{\|\mathbf{X} - \mathbf{Z}\hat{\theta}(\mathbf{X})\|^2 k} < f_\gamma \right\}.$$

Стоит заметить, что геометрически доверительная область будет иметь форму эллипсоида в \mathbb{R}^k .

Упражнение

Представить выборку из нормального распределения $\mathcal{N}(a, \sigma^2)$ в виде гауссовской линейной модели и построить доверительные интервалы для параметров с помощью общей теории.

15.

Лекция 15

На этой лекции мы начнем разбирать последнюю большую тему нашего курса - проверку статистических гипотез. Разберемся сначала с постановкой проблемы.

Пусть \mathbf{X} - наблюдение с неизвестным распределением $P \in \mathcal{P}$, где \mathcal{P} есть некоторое семейство распределений.

Далее, возьмем 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 дальнейший анализ.

Если ответ положительный, то \mathcal{P} можно сузить до \mathcal{P}_0 . Если же ответ отрицательный, то \mathcal{P} сужается до $\mathcal{P} \setminus \mathcal{P}_0$. В общем случае может быть выдвинуто сразу несколько гипотез, из которых нам необходимо выбрать одну, чтобы "локализовать" истинное распределение.

Гипотезы

Определение 15.1. Статистической гипотезой называется предположение вида

$$H_0 : P \in \mathcal{P}_0,$$

где $\mathcal{P}_0 \subset \mathcal{P}$ - подсемейство распределений.

Задача

По наблюдению \mathbf{X} либо принять гипотезу (тогда \mathcal{P} заменяется на \mathcal{P}_0), либо отклонить (тогда \mathcal{P} заменяется на $\mathcal{P} \setminus \mathcal{P}_0$).

Во втором случае мы можем выдвинуть альтернативную гипотезу $H_1 : P \in \mathcal{P}_1$, где $\mathcal{P}_1 \subseteq \mathcal{P} \setminus \mathcal{P}_0$. Соответственно, если у нас есть альтернатива, то мы должны проверять гипотезу H_0 против альтернативы H_1 . Альтернатив может быть несколько:

$$H_0 : P \in \mathcal{P}_0,$$

$$H_1 : P \in \mathcal{P}_1,$$

$$\vdots$$

$$H_k : P \in \mathcal{P}_k.$$

Статистические критерии

Методология проверки гипотез следующая.

Определение 15.2. Пусть \mathcal{X} - выборочное пространство, то есть множество всех возможных значений наблюдений. Тогда подмножество $S \subset \mathcal{X}$ называется критерием или критическим множеством для проверки гипотезы H_0 (против альтернативы H_1), если правило принятия H_0 выглядит так:

H_0 отвергается тогда и только тогда, когда $\mathbf{X} \in S$.

Тем самым, S является "плохим" множеством, и, если \mathbf{X} попало в него, то H_0 признается неверной и выбирается альтернатива.

Замечание

Множество $\mathcal{X} \setminus S$ принято называть областью принятия гипотезы.

Естественно, что при принятии решений возможны ошибки. Например, мы можем промахнуться следующим образом: гипотезу H_0 приняли, хотя, на самом деле, она была не верна. Принято различать ошибки первого и второго рода.

Ошибки первого и второго рода

Определение 15.3. Пусть проверяется гипотеза H_0 . 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 неверна.

Вопрос: какая из этих ошибок опаснее?

Методологически считается, что ошибка первого рода опаснее.

Идея состоит в том, что мы считаем, что при выдвижении гипотезы H_0 мы выбрали наше подсемейство \mathcal{P}_0 из тех соображений, что мы не хотим совершать ошибку первого рода. Предположим, что мы неверно отвергли H_0 . Тогда мы вынуждены рассматривать некоторую альтернативу H_1 , которая заведомо неверна. Последнее приведет к тому, что мы уйдем далеко от правильного ответа, будет проводиться много бесполезной работы.

В случае совершения ошибки второго рода мы принимаем H_0 , но она неверна. В таком случае вердикт может измениться, если добавятся новые данные, и мы сможем двинуться в сторону чего-то более правильного. Другими словами, здесь шансы совершить значительное количество лишней работы заметно меньше.

На этой методологии основан принцип сравнения критериев.

Функция мощности

Задача

Нужно предложить статистический критерий с небольшими вероятностями ошибок как первого, так и второго рода.

Определение 15.4. Пусть S - критерий для проверки гипотезы H_0 . Функцией мощности критерия S называется

$$\beta(Q, S) = Q(\mathbf{X} \in S), \quad Q \in \mathcal{P}.$$

Возникает интересный момент: вероятность ошибки не является одной вероятностью, а является набором вероятностей.

Наблюдение

Вероятностями ошибок первого рода называется набор $\{\beta(Q, S) : Q \in \mathcal{P}_0\}$.

Вероятностями ошибок второго рода называется набор $\{1 - \beta(Q, S) : Q \in \mathcal{P}_1\}$ (или $Q \in \mathcal{P} \setminus \mathcal{P}_0$).

Соответственно, здесь мы получаем следующую постановку: мы хотим критерий, который имел бы маленькие значения функции мощности на основной гипотезе и большие на альтернативе. В идеале, хочется получить нули, но эта задача безнадежна, и нужно на что-то соглашаться. Но на что? И как определить, что один критерий лучше другого? Обсудим это.

Уровень значимости

Разумно рассматривать только критерии, у которых вероятности ошибки первого рода равномерно ограничены сверху каким-то ε , раз уж мы решили, что ошибка первого рода опаснее.

Определение 15.5. Пусть $S \subset \mathcal{X}$ - это критерий для проверки гипотезы $H_0 : P \in \mathcal{P}_0$. Будем говорить, что критерий S имеет уровень значимости ε , если для любого $Q \in \mathcal{P}_0$ выполнено

$$\beta(Q, S) \leq \varepsilon.$$

По сути, если ε является подходящим уровнем значимости, то и любое число, большее ε , тоже подходит. Поэтому, наряду с уровнем значимости вводится минимальный уровень значимости:

Определение 15.6. Минимальным уровнем значимости или размером критерия называют

$$\alpha(S) = \sup_{Q \in \mathcal{P}_0} \beta(Q, S).$$

Несмещенность и состоятельность оценки

Опираясь на функцию мощности можно ввести два весьма полезных свойства.

Определение 15.7. Критерий S для проверки гипотезы $H_0 : P \in \mathcal{P}_0$ против альтернативы $H_1 : P \in \mathcal{P}_1$ называется несмещенным, если

$$\sup_{Q \in \mathcal{P}_0} \beta(Q, S) \leq \inf_{Q \in \mathcal{P}_1} \beta(Q, S).$$

Данное свойство достаточно естественно: мы хотим, чтобы на \mathcal{P}_0 значение функции мощности было маленьким, а на \mathcal{P}_1 большим.

Определение 15.8. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка растущего размера. Тогда последовательность критериев $\{S_n, n \in \mathbb{N}\}$, $S_n \subset \mathcal{X}^n$, для проверки гипотезы $H_0 : P \in \mathcal{P}_0$ против альтернативы $H_1 : P \in \mathcal{P}_1$ называется состоятельной, если для любого $Q \in \mathcal{P}_1$ выполнено

$$\beta(Q, S_n) \rightarrow 1 \text{ при } n \rightarrow \infty.$$

Смысл: вероятности ошибки второго рода стремятся к нулю.

Сравнение критериев

Обсудим теперь сравнение критериев. Пусть есть два критерия S и R . Как можно сравнить их и понять, какой лучше?

Важный принцип. Мы должны выбрать некоторый уровень значимости ε , вероятность ошибки первого рода, на которую мы согласились, и сравнивать только критерии этого уровня значимости.

Замечание

Однако, "загонять" ε в ноль неосмысленно, так как это лишь увеличит вероятность ошибки второго рода. Здесь нужно искать баланс.

Определение 15.9. Пусть S и R - критерии уровня значимости ε для проверки гипотезы $H_0 : P \in \mathcal{P}_0$ против альтернативы $H_1 : P \in \mathcal{P}_1$. Будем говорить, что S (равномерно) мощнее R , если для любого $Q \in \mathcal{P}_1$ выполнено

$$\beta(Q, S) \geq \beta(Q, R),$$

то есть вероятность ошибки второго рода у критерия S равномерно меньше, чем у R .

Равномерно наиболее мощные критерии

Если же критерий оказывается мощнее всех остальных, то мы его будем называть равномерно наиболее мощным.

Определение 15.10. Критерий S для проверки гипотезы $H_0 : P \in \mathcal{P}_0$ против альтернативы $H_1 P \in \mathcal{P}_1$ называется равномерно наиболее мощным критерием (р.н.м.к.) уровня значимости ε , если $\alpha(S) \leq \varepsilon$ и S мощнее любого другого критерия R уровня значимости ε .

Вопрос: как строить р.н.м.к.?

Оказывается в случае простых гипотез есть явный ответ.

Определение 15.11. Гипотеза вида $H_0 : P = P_0$, где P_0 - некоторое известное распределение, называется простой.

Лемма Неймана-Пирсона

Рассмотрим проверку простой гипотезы $H_0 : P = P_0$ против простой альтернативы $H_1 : P = P_1$. 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 должны быть взяты из одного доминируемого семейства.

Замечание

Например, различить пуассоновское и экспоненциальное распределения крайне просто - считаем, что экспоненциальное, если получили нецелое число, и пуассоновское - иначе. Такой критерий никогда не ошибается.

Для $\lambda > 0$ рассмотрим критерий S_λ :

$$S_\lambda = \{x \in \mathcal{X} : p_1(x) - \lambda p_0(x) \geq 0\}.$$

Оказывается, такой критерий обладает интересными свойствами.

Лемма (Нейман-Пирсон)

Пусть критерий R таков, что $P_0(\mathbf{X} \in R) \leq P_0(\mathbf{X} \in S_\lambda)$. Тогда

- 1) $P_1(\mathbf{X} \in R) \leq P_1(\mathbf{X} \in S_\lambda)$,
- 2) $P_0(\mathbf{X} \in S_\lambda) \leq P_1(\mathbf{X} \in S_\lambda)$.

Доказательство.

1) Для начала заметим, что для любого $x \in \mathcal{X}$

$$I_R(x)(p_1(x) - \lambda p_0(x)) \leq I_{S_\lambda}(x)(p_1(x) - \lambda p_0(x)).$$

Это верно для всех x : предположим, что x попадает в S_λ , тогда $I_{S_\lambda} = 1$, $I_R \leq 1$, а $p_1(x) - \lambda p_0(x) \geq 0$. Если же x не попадает в S_λ , то $I_{S_\lambda} = 0$, $I_R \geq 0$, а $p_1(x) - \lambda p_0(x) \leq 0$.

Проинтегрируем неравенства по мере μ по всему выборочному пространству:

$$\int_R (p_1(x) - \lambda p_0(x)) \mu(dx) \leq \int_{S_\lambda} (p_1(x) - \lambda p_0(x)) \mu(dx).$$

Раскроем интегралы по линейности, вынесем константы, тогда:

$$P_1(\mathbf{X} \in R) - \lambda \cdot P_0(\mathbf{X} \in R) \leq P_1(\mathbf{X} \in S_\lambda) - \lambda \cdot P_0(\mathbf{X} \in S_\lambda).$$

Следовательно,

$$P_1(\mathbf{X} \in R) - P_1(\mathbf{X} \in S_\lambda) \leq \lambda \cdot (P_0(\mathbf{X} \in R) - P_0(\mathbf{X} \in S_\lambda)) \leq 0.$$

Первое неравенство доказано.

2) Проверим, что для любого $\lambda > 0$ выполнено $P_0(\mathbf{X} \in S_\lambda) \leq P_1(\mathbf{X} \in S_\lambda)$.

Рассмотрим два случая: $\lambda \geq 1$ и $\lambda < 1$.

- 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,

$$P_0(\mathbf{X} \in S_\lambda) = \int_{S_\lambda} p_0(x) \mu(dx) \leq \int_{S_\lambda} p_1(x) \mu(dx) = P_1(\mathbf{X} \in S_\lambda).$$

- Пусть $\lambda \in (0, 1)$. Тогда для $x \notin S_\lambda$ выполнено

$$p_1(x) < \lambda p_0(x) < p_0(x).$$

Следовательно,

$$P_1(\mathbf{X} \in \overline{S_\lambda}) = \int_{\overline{S_\lambda}} p_1(x) \mu(dx) \leq \int_{\overline{S_\lambda}} p_0(x) \mu(dx) = P_0(\mathbf{X} \in \overline{S_\lambda}),$$

что эквивалентно $P_0(\mathbf{X} \in S_\lambda) \leq P_1(\mathbf{X} \in S_\lambda)$.

Следствия

У данной леммы есть два важных следствия. Первое следствие говорит нам о том, как мы можем построить равномерно наиболее мощный критерий уровня значимости ε .

Следствие (1)

Если $\lambda > 0$ удовлетворяет соотношению $P_0(\mathbf{X} \in S_\lambda) = \varepsilon$, то S_λ является равномерно наиболее мощным критерием уровня значимости ε для проверки гипотезы $H_0 : P = P_0$ против альтернативы $H_1 : P = P_1$.

Доказательство.

Если R - любой другой критерий уровня значимости $\varepsilon > 0$, то $P_0(\mathbf{X} \in R) \leq \varepsilon = P_0(\mathbf{X} \in S_\lambda)$. По первому утверждению леммы Неймана-Пирсона получаем, что $P_1(\mathbf{X} \in R) \leq P_1(\mathbf{X} \in S_\lambda)$, то есть S_λ мощнее R . ■

Следствие (2)

Критерий S_λ несмещенный.

Доказательство. Это второе утверждение леммы Неймана-Пирсона. ■

Тем самым мы получили способ, с помощью которого можно построить р.н.м.к. для проверки простой гипотезы против простой альтернативы.

Из следствия (1) можно сделать вывод о том, что для нахождения р.н.м.к. необходимо решить (относительно λ) уравнение

$$P_0(\mathbf{X} \in S_\lambda) = \int_{S_\lambda} p_0(x) \mu(Dx) = \varepsilon,$$

где $S_\lambda = \{x \in \mathcal{X} : p_1(x) - \lambda p_0(x) \geq 0\}$. Оно почти всегда разрешимо в случае абсолютно непрерывных распределений. В дискретном случае, наоборот, оно для почти всех ε неразрешимо. В этом случае стоит изменить ε так, чтобы оно стало разрешимо.

Комментарий

В общей ситуации можно ввести рандомизированные критерии. Мы не будем подробно на них останавливаться, лишь кратко опишем идею. Для каждого значения наблюдения мы выбираем с какой вероятностью мы примем гипотезу H_0 против альтернативы H_1 , то есть получается некоторая функция от x , которая принимает

значения на $[0, 1]$. Если значение не 0 или 1, то мы не можем однозначно выбрать гипотезу. Тогда мы "подбрасываем монетку" с вероятностью, равной значению функции.

Польза от такого подхода следующая. Например, у нас дискретное распределение, следовательно вероятность P_0 как-то дискретно меняется, и для некоторого λ вероятность меньше ε , а для другого λ больше ε , и эти два значения отличаются на событие, что x попал в какие-то дискретные точки. Соответственно, для $\lambda < \varepsilon$ мы будем предпочитать гипотезу H_1 , для всех, которые не попали во второй случай, будем предпочитать H_0 , а для тех, которые оказались в промежутке, будем пытаться подобрать числа так, чтобы вероятность P_0 была бы в точности равна ε . Это станет более понятно, когда мы дойдем до примеров.

Монотонное отношение правдоподобия

Мы построили р.н.м.к. для проверки простой гипотезы против простой альтернативы. Возникает естественный вопрос: можно ли обобщить данный результат на случай сложных гипотез?

Оказывается, это можно сделать, когда речь идет о параметрических семействах с монотонным отношением правдоподобия. В параметрической модели $\{P_\theta, \theta \in \Theta\}$ гипотезы естественно записывать, как предположения о значении параметра. Например,

$$H_0 : \theta \in \Theta_0.$$

Функцию мощности критерия S также удобно записывать, как функцию от параметра $\beta(\theta, S) = P_\theta(\mathbf{X} \in S)$.

Итак, пусть \mathbf{X} - это наблюдение с неизвестным распределением P , принадлежащим параметрическому семейству $\{P_\theta, \theta \in \Theta\}$, причем семейство $\{P_\theta, \theta \in \Theta\}$ является доминируемым с плотностью $p_\theta(\mathbf{x})$.

Определение 15.12. Пусть $\{P_\theta, \theta \in \Theta\}$, $\Theta \subseteq \mathbb{R}$ - это доминируемое семейство распределений с плотностью $p_\theta(x)$ по мере μ . Тогда это семейство имеет монотонное отношение правдоподобия по статистике $T(\mathbf{X}) \in \mathbb{R}$, если для любых $\theta_1 < \theta_2$, $\theta_1, \theta_2 \in \Theta$, отношение плотностей имеет вид

$$\frac{p_{\theta_2}(\mathbf{X})}{p_{\theta_1}(\mathbf{X})} = \psi_{\theta_1, \theta_2}(T(\mathbf{X})), \quad (1)$$

где $\psi_{\theta_1, \theta_2}(t)$ - это монотонная функция, причем ее монотонность одинакова для всех $\theta_1 < \theta_2$.

Замечание

В этом случае $T(\mathbf{X})$ - это достаточная статистика.

Лемма

Пусть $\psi_{\theta_1, \theta_2}(t)$ в (1) всегда не убывает по $T(\mathbf{X})$. Тогда для всех $c \in \mathbb{R}$ и для всех $\theta_1 < \theta_2$, $\theta_1, \theta_2 \in \Theta$ выполняется неравенство

$$P_{\theta_1}(T(\mathbf{X}) \geq c) \leq P_{\theta_2}(T(\mathbf{X}) \geq c),$$

то есть функция $P_{\theta}(T(\mathbf{X}) \geq c)$ не убывает по θ .

Доказательство.

Считаем, что $\psi_{\theta_1, \theta_2}(t) \geq 0$ и не убывает на \mathbb{R} . Введем множество $D = \{\mathbf{x} : T(\mathbf{x}) \geq c\}$. Если $\psi_{\theta_1, \theta_2}(c) \geq 1$, то для любого $\mathbf{x} \in D$ выполнено

$$p_{\theta_2}(\mathbf{x}) = \psi_{\theta_1, \theta_2}(T(\mathbf{x})) \cdot p_{\theta_1}(\mathbf{x}) \geq \psi_{\theta_1, \theta_2}(c) \cdot p_{\theta_1}(\mathbf{x}) \geq p_{\theta_1}(\mathbf{x}).$$

Интегрируем это неравенство по мере μ по множеству D :

$$P_{\theta_2}(T(\mathbf{X}) \geq c) = \int_D p_{\theta_2}(\mathbf{x}) \mu(d\mathbf{x}) \geq \int_D p_{\theta_1}(\mathbf{x}) \mu(d\mathbf{x}) = P_{\theta_1}(T(\mathbf{X}) \geq c).$$

Если же c таково, что $\psi_{\theta_1, \theta_2}(c) \in [0, 1]$, то рассматриваем множество \bar{D} . На нем, наоборот, $p_{\theta_2}(\mathbf{x}) \leq p_{\theta_1}(\mathbf{x})$ и, стало быть,

$$P_{\theta_2}(T(\mathbf{X}) < c) = \int_{\bar{D}} p_{\theta_2}(\mathbf{x}) \mu(d\mathbf{x}) \leq \int_{\bar{D}} p_{\theta_1}(\mathbf{x}) \mu(d\mathbf{x}) = P_{\theta_1}(T(\mathbf{X}) < c).$$

В итоге, получаем искомое:

$$P_{\theta_1}(T(\mathbf{X}) \geq c) \leq P_{\theta_2}(T(\mathbf{X}) \geq c).$$

■

Эта лемма лежит в основе доказательства следующей теоремы, которая позволяет нам строить р.н.м.к. в ситуации, когда мы имеем дело с параметрическими семействами с монотонным отношением правдоподобия, и в ситуации, когда мы рассматриваем "односторонние" гипотезы.

Теорема 15.1. (о монотонном отношении правдоподобия)

Пусть $\{P_{\theta}, \theta \in \Theta\}$, $\theta \in \mathbb{R}$ - семейство с монотонным отношением правдоподобия по

статистике $T(\mathbf{X})$, причем функция $\psi_{\theta_1, \theta_2}(t)$ в (1) не убывает и непрерывна. Пусть $\theta_0 \in \Theta$. Если $c \in \mathbb{R}$ удовлетворяет соотношению

$$P_{\theta_0}(T(\mathbf{X}) \geq c) = \varepsilon,$$

то критерий $S = \{\mathbf{x} : T(\mathbf{x}) \geq c\}$ является равномерно наиболее мощным критерием уровня значимости ε для проверки гипотезы $H_0 : \theta \leq \theta_0$ против альтернативы $H_1 : \theta > \theta_0$.

Доказательство.

Для начала возьмем какое-нибудь $\theta < \theta_0$. Тогда по лемме

$$P_{\theta}(T(\mathbf{X}) \geq c) \leq P_{\theta_0}(T(\mathbf{X}) \geq c) = \varepsilon.$$

Из этого следует, что критерий S имеет уровень значимости ε . Проверим теперь, что он является равномерно наиболее мощным.

Пусть R - любой другой критерий с уровнем значимости ε . Возьмем любое $\theta_1 > \theta_0$. Нужно доказать, что

$$P_{\theta_1}(\mathbf{X} \in R) \leq P_{\theta_1}(\mathbf{X} \in S).$$

Согласно лемме Неймана-Пирсона равномерно наиболее мощный критерий уровня значимости ε для проверки гипотезы $H_0 : \theta = \theta_0$ против гипотезы $H_1 : \theta = \theta_1$ имеет вид

$$S_{\lambda} = \{\mathbf{x} : p_{\theta_1}(\mathbf{x}) - \lambda p_{\theta_0}(\mathbf{x}) \geq 0\}, \text{ где } P_{\theta_0}(\mathbf{X} \in S_{\lambda}) = \varepsilon.$$

Однако в силу непрерывности и неубывания $\psi_{\theta_0, \theta_1}$ для некоторого $\tilde{\lambda}$ выполнено

$$S_{\lambda} = \left\{ \mathbf{x} : \frac{p_{\theta_1}(\mathbf{x})}{p_{\theta_0}(\mathbf{x})} \geq \lambda \right\} = \{\mathbf{x} : \psi_{\theta_0, \theta_1}(T(\mathbf{x})) \geq \lambda\} = \{\mathbf{x} : T(\mathbf{x}) \geq \tilde{\lambda}\}.$$

Заметим, что $\tilde{\lambda} = e$ подходит под условие. Следовательно, для любого $\theta_1 > \theta_0$

$$P_{\theta_1}(\mathbf{X} \in R) \leq P_{\theta_1}(\mathbf{X} \in S_{\lambda}) = P_{\theta_1}(T(\mathbf{X}) \geq e) = P_{\theta_1}(\mathbf{X} \in S).$$

Значит, S является равномерно наиболее мощным критерием для проверки гипотезы $H_0 : \theta \leq \theta_0$ против альтернативы $H_1 : \theta > \theta_0$. ■

Замечание (1)

Тот же самый критерий S является равномерно наиболее мощным критерием уровня значимости ε для проверки гипотезы $H_0 : \theta = \theta_0$ против альтернативы $H_1 : \theta > \theta_0$.

Замечание (2)

Теорема означает невозможность построения р.н.м.к. в случае двухсторонней альтернативы, то есть для $H_0 : \theta = \theta_0$ против $H_1 : \theta \neq \theta_0$, когда мы имеем дело с

экспоненциальными семействами.

Замечание (3)

В случае убывания функции $\psi_{\theta_1, \theta_2}$ р.н.м.к. будет иметь вид $T(\mathbf{X}) \leq e$. Он же будет оптимальным, если основная гипотеза и альтернатива поменяются местами: $H_0 : \theta \geq \theta_0$ против альтернативы $H_1 : \theta < \theta_0$. В общем случае можно просто решать неравенство

$$p_{\theta_1}(x) - \lambda p_{\theta_0}(x) \geq 0,$$

где θ_1 берется из альтернативы, а θ_0 - из основной гипотезы.

На следующей лекции мы разберемся с тем, как можно проверять линейные гипотезы в линейной гауссовской модели, и разберем некоторые примеры.

16.

Лекция 16

Примеры применения теоремы о монотонном отношении правдоподобия

Начнем эту лекцию с примеров, которые не были разобраны на прошлой лекции. Также обсудим что можно сделать в ситуации, когда мы не можем честно решить уравнение на λ , которое возникает при применении монотонного отношения правдоподобия.

Пример (1)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из схемы Бернулли $\text{Bin}(1, \theta)$. Постройте р.н.м.к. уровня значимости ε для проверки гипотезы $H_0 : \theta \leq \theta_0$ против альтернативы $H_1 : \theta > \theta_0$. Проверьте этот критерий на состоятельность.

Решение

Выпишем функцию правдоподобия:

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

Согласно теореме о монотонном отношении правдоподобия необходимо найти отношение правдоподобия: при $\theta_1 > \theta_0$

$$\frac{p_{\theta_1}(\mathbf{X})}{p_{\theta_0}(\mathbf{X})} = \frac{\theta_1^{n\bar{X}} (1 - \theta_1)^{n - n\bar{X}}}{\theta_0^{n\bar{X}} (1 - \theta_0)^{n - n\bar{X}}} = \left(\frac{\theta_1 (1 - \theta_0)}{\theta_0 (1 - \theta_1)} \right)^{n\bar{X}} \left(\frac{1 - \theta_1}{1 - \theta_0} \right)^n.$$

Решаем неравенство:

$$\frac{p_{\theta_1}(\mathbf{X})}{p_{\theta_0}(\mathbf{X})} = \left(\frac{\theta_1 (1 - \theta_0)}{\theta_0 (1 - \theta_1)} \right)^{n\bar{X}} \left(\frac{1 - \theta_1}{1 - \theta_0} \right)^n \geq \lambda.$$

Так как $\theta_1 > \theta_0$, то выражение в левой части возрастает по $\sum_{i=1}^n X_i$. Тогда неравенство выше эквивалентно тому, что

$$\sum_{i=1}^n X_i \geq c_1 = c_1(\lambda, \theta_0, \theta_1).$$

Теперь надо найти c_1 из уравнения

$$P_{\theta_0} \left(\sum_{i=1}^n X_i \geq c_1 \right) = \varepsilon.$$

В силу того, что $\sum_{i=1}^n X_i \sim \text{Bin}(n, \theta_0)$, находим, что c_1 - это $(1 - \varepsilon)$ -квантиль распределения $\text{Bin}(n, \theta_0)$.

Мы нашли вид критерия и чему равняется c_1 . Проверим теперь состоятельность критерия. Пусть $\theta_1 > \theta_0$. Нам нужно проверить, что

$$P_{\theta_1} \left(\sum_{i=1}^n X_i \geq c_1 \right) \rightarrow 1.$$

Мы знаем, что если θ_1 - это истинное значение параметра, то $\sum_{i=1}^n X_i \sim \text{Bin}(n, \theta_1)$, но мы не понимаем, чему равно c_1 . Перепишем условие на вероятность ошибки первого рода:

$$\varepsilon = P_{\theta_0} \left(\sum_{i=1}^n X_i \geq c_1 \right) = P_{\theta_0} \left(\sqrt{n} \frac{(\bar{X} - \theta_0)}{\sqrt{\theta_0(1 - \theta_0)}} \geq \sqrt{n} \frac{\frac{c_1}{n - \theta_0}}{\sqrt{\theta_0(1 - \theta_0)}} \right).$$

Левая часть в неравенстве по ЦПТ сходится по распределению к стандартному нормальному закону. Стало быть,

$$\sqrt{n} \frac{\frac{c_1}{n - \theta_0}}{\sqrt{\theta_0(1 - \theta_0)}} \rightarrow u_{1-\varepsilon}$$

при $n \rightarrow +\infty$, где $u_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль $\mathcal{N}(0, 1)$. Отсюда получаем, что $c_1 = n\theta_0 + O(\sqrt{n})$. Следовательно,

$$P_{\theta_1} \left(\sum_{i=1}^n X_i \geq c_1 \right) = P_{\theta_1} \left(\sqrt{n} \frac{(\bar{X} - \theta_1)}{\sqrt{\theta_1(1 - \theta_1)}} \geq \sqrt{n} \frac{\frac{c_1}{n - \theta_1}}{\sqrt{\theta_1(1 - \theta_1)}} \right) \rightarrow 1,$$

так как левая часть снова сходится по распределению к $\mathcal{N}(0, 1)$, а правая часть стремится к $-\infty$ в силу условия $\theta_1 > \theta_0$. Следовательно, критерий состоятелен.

На примере биномиального распределения удобно посмотреть на то, как бы работал рандомизированный критерий проверки гипотез.

У нас были гипотезы $H_0 : \theta \leq \theta_0$ и $H_1 : \theta > \theta_0$, и мы получили критерий $\sum_{i=1}^n X_i \geq c_1$, где c_1 находится из соображений, что $P_{\theta_0} \left(\sum_{i=1}^n X_i \geq c_1 \right) = \varepsilon$.

Далее, если ε таково, что ни для какого c_1 мы не можем получить точное равенство, то нам следует уменьшить ε , чтобы уравнение стало разрешимым. Но предположим, что мы хотим решить задачу для всех ε . Тогда стоит действовать так: пусть $c \in \mathbb{N}$ такое, что

$$P_{\theta_0} \left(\sum X_i \geq c + 1 \right) < \varepsilon < P_{\theta_0} \left(\sum X_i \geq c \right).$$

Тогда введем функцию $\psi(X_1, \dots, X_n) = \begin{cases} 1, & \sum X_i \geq c + 1 \\ 0, & \sum X_i < c \\ \rho, & \sum X_i = c \end{cases}$. ρ выбирается таким образом, что $E_{\theta_0} \psi(X_1, \dots, X_n) = \varepsilon$, то есть ρ находится из уравнения

$$P_{\theta_0} \left(\sum_{i=1}^n X_i \geq c + 1 \right) + \rho P_{\theta_0} \left(\sum_{i=1}^n X_i = c \right) = \varepsilon$$

Тогда при $\rho = 1$ левая часть становится больше ε , а при $\rho = 0$ - меньше ε , поэтому найдется $\rho \in (0, 1)$, для которого равенство выполняется. Соответственно, правило принятия гипотезы говорит о том, что если $\psi = 1$, то мы выбираем гипотезу H_1 , если $\psi = 0$ - гипотезу H_0 , а если что-то между 0 и 1, то мы "бросаем монетку", которая с вероятностью ρ предпочитает H_1 , а с вероятностью $1 - \rho$ предпочитает H_0 . Продолжим разбирать примеры.

Пример (2)

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из нормального распределения $\mathcal{N}(0, 1)$. Постройте р.н.м.к. уровня значимости ε для проверки гипотезы $H_0 : \theta \geq \theta_0$ против альтернативы $H_1 : \theta < \theta_0$.

Решение

Выпишем функцию правдоподобия:

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

Необходимо найти отношение правдоподобия, где $\theta_1 < \theta_0$:

$$\frac{p_{\theta_1}(\mathbf{X})}{p_{\theta_0}(\mathbf{X})} = \frac{\exp \left\{ -\sum_{i=1}^n \frac{(X_i - \theta_1)^2}{2} \right\}}{\exp \left\{ -\sum_{i=1}^n \frac{(X_i - \theta_0)^2}{2} \right\}} = \exp \left\{ (\theta_1 - \theta_0) \sum_{i=1}^n X_i + \frac{n}{2}(\theta_0^2 - \theta_1^2) \right\}.$$

Решаем неравенство:

$$\frac{p_{\theta_1}(\mathbf{X})}{p_{\theta_0}(\mathbf{X})} = \exp \left\{ (\theta_1 - \theta_0) \sum_{i=1}^n X_i + \frac{n}{2}(\theta_0^2 - \theta_1^2) \right\} \geq \lambda.$$

Раз $\theta_1 < \theta_0$, то выражение в левой части убывает по $\sum_{i=1}^n X_i$. Тогда неравенство выше эквивалентно тому, что $\sum_{i=1}^n X_i \leq c = c(\lambda, \theta_0, \theta_1)$. Теперь надо найти c из уравнения

$$P_{\theta_0} \left(\sum_{i=1}^n X_i \leq c \right) = \varepsilon.$$

В случае нормального распределения всегда удобно переходить к стандартному нормальному распределению. Заметим, что $\sqrt{n}(\bar{\mathbf{X}} - \theta_0) \sim \mathcal{N}(0, 1)$, если θ_0 - истинное значение. Отсюда, критерий имеет вид

$$S = \{\sqrt{n}(\bar{\mathbf{X}} - \theta_0) \leq u\},$$

где u - это ε -квантиль распределения $\mathcal{N}(0, 1)$.

В этом примере видно, что для любого ε подобное u находится без проблем, так как распределение абсолютно непрерывно.

Двойственность проверки гипотез и доверительного оценивания

Прежде чем переходить к обсуждению линейных гипотез в линейной гауссовской модели, поговорим о двойственности проверки гипотез и доверительного оценивания. Идея проверить гипотезу и идея оценить параметр, особенно если гипотеза записывается в виде предположения о значении параметра, выглядят очень близкими задачами. Умение строить доверительные интервалы помогает при проверке гипотез, и наоборот. Пусть \mathbf{X} - это наблюдение с неизвестным распределением P , принадлежащим параметрическому семейству $\{P_\theta, \theta \in \Theta\}$.

1. Пусть $S(\mathbf{X}) \subset \Theta$ - доверительная область уровня доверия $\gamma = 1 - \varepsilon$ для параметра θ . Хотим проверить простую гипотезу $H_0 : \theta = \theta_0$.

Утверждение (1)

Критерий $\tilde{S}_{\theta_0} = \{\mathbf{X} : \theta_0 \notin S(\mathbf{X})\}$ является критерием уровня значимости ε для проверки гипотезы H_0 .

Доказательство.

$$P_{\theta_0}(\mathbf{X} \in \tilde{S}_{\theta_0}) = P_{\theta_0}(\theta_0 \notin S(\mathbf{X})) = 1 - P_{\theta_0}(\theta_0 \in S(\mathbf{X})) \leq \varepsilon.$$

■

2. Пусть для любого $\theta_0 \in \Theta$ мы умеем строить критерий \tilde{S}_{θ_0} уровня значимости ε для проверки простой гипотезы $H_0 : \theta = \theta_0$. Хотим построить доверительную область.

Утверждение (2)

Область $S(\mathbf{X}) = \{\theta \in \Theta : \mathbf{X} \notin \tilde{S}_\theta\}$ является доверительной областью уровня доверия $\gamma = 1 - \varepsilon$ для параметра θ .

Доказательство. Для любого $\theta_0 \in \Theta$ выполнено

$$P_{\theta_0}(\theta_0 \in S(\mathbf{X})) = P_{\theta_0}(\mathbf{X} \notin \tilde{S}_{\theta_0}) = 1 - P_{\theta_0}(\mathbf{X} \in \tilde{S}_{\theta_0}) \geq 1 - \varepsilon.$$

■

Применим эти соображения при проверке линейных гипотез в линейной гауссовской регрессии.

Линейные гипотезы

Пусть

$$\mathbf{X} = \mathbf{Z}\boldsymbol{\theta} + \varepsilon$$

- гауссовская линейная модель $\mathbf{Z} \in \text{Mat}(n \times k)$, $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Неизвестные параметры: $\boldsymbol{\theta} \in \mathbb{R}^k$, $\sigma > 0$.

Задача

Построить критерий для проверки линейной гипотезы:

$$H_0 : \mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau},$$

где $\mathbf{A} = \text{Mat}(m \times k)$, $\boldsymbol{\tau} \in \mathbb{R}^m$ - известны и $\text{rank } \mathbf{A} = m \leq k$.

Решение

Мы хотим проверить гипотезу следующего вида: изначально $\boldsymbol{\theta} \in \mathbb{R}^k$, но мы хотим проверить, что $\boldsymbol{\theta}$ пробегает лишь некоторое аффинное подпространство, которое задается уравнением $\mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau}$.

Как строить? Рассмотрим оценку наименьших квадратов $\hat{\boldsymbol{\theta}}(\mathbf{X}) = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X}$. Тогда мы знаем, что $\hat{\boldsymbol{\tau}} = \hat{\boldsymbol{\tau}}(\mathbf{X}) = \mathbf{A}\hat{\boldsymbol{\theta}}(\mathbf{X})$ является оптимальной оценкой $\mathbf{A}\boldsymbol{\theta}$.

Лемма

Обозначим $\mathbf{D} = \mathbf{A}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{A}^T$ и положим

$$Q_A(\mathbf{X}) = (\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta})^T \mathbf{D}^{-1} (\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta}).$$

Тогда

$$\frac{Q_A(\mathbf{X})}{\sigma^2} \sim \chi_m^2.$$

Доказательство.

Напомним, что $\hat{\boldsymbol{\theta}}(\mathbf{X}) \sim \mathcal{N}(\boldsymbol{\theta}, \sigma^2(\mathbf{Z}^T \mathbf{Z})^{-1})$. Тогда, конечно,

$$\hat{\boldsymbol{\tau}} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\theta}, \sigma^2 \mathbf{A}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{A}^T) = \mathcal{N}(\mathbf{A}\boldsymbol{\theta}, \sigma^2 \mathbf{D}).$$

Матрица $\mathbf{D} \in \text{Mat}(m \times m)$ положительно определена. Значит, из нее можно извлечь корень: $\mathbf{D}^{\frac{1}{2}}$, где $(\mathbf{D}^{\frac{1}{2}})^2 = \mathbf{D}$ и $\mathbf{D}^{\frac{1}{2}} = (\mathbf{D}^{\frac{1}{2}})^T$. Тогда

$$\frac{1}{\sigma} \cdot \mathbf{D}^{\frac{1}{2}}(\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta}) \sim \mathcal{N}(\mathbf{0}, \mathbf{D}^{-\frac{1}{2}} \cdot \mathbf{D} \cdot \mathbf{D}^{-\frac{1}{2}}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_m).$$

В итоге,

$$\left| \frac{1}{\sigma} \mathbf{D}^{-\frac{1}{2}}(\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta}) \right| \sim \chi_m^2.$$

Осталось заметить, что

$$\left| \mathbf{D}^{-\frac{1}{2}}(\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta}) \right|^2 = (\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta})^T \mathbf{D}^{-\frac{1}{2}} \cdot \mathbf{D}^{-\frac{1}{2}}(\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta}) = (\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta})^T \mathbf{D}^{-1}(\hat{\boldsymbol{\tau}} - \mathbf{A}\boldsymbol{\theta}) = Q_A(\mathbf{X}).$$

■

Построение критерия

Рассмотрим статистику, похожую на $Q_A(\mathbf{X})$:

$$\hat{Q}_A(\mathbf{X}) = (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau})^T \mathbf{D}^{-1}(\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}).$$

Согласно лемме в условиях гипотезы $H_0 : \mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau}$ выполнено

$$\frac{1}{\sigma^2} \hat{Q}_A(\mathbf{X}) \sim \chi_m^2.$$

Также заметим, что статистика $\hat{Q}_A(\mathbf{X})$ является функцией от $\hat{\boldsymbol{\theta}}(\mathbf{X})$. Следовательно, она независима со статистикой $\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2$.

Далее, вспомним, что $\frac{1}{\sigma^2} \|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 \sim \chi_{n-k}^2$. Стало быть, если гипотеза H_0 верна, то выполнено

$$\frac{\hat{Q}_A(\mathbf{X})}{\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2} \cdot \frac{n-k}{m} \sim F_{m, n-k}.$$

Теперь мы готовы сформулировать F-критерий для проверки гипотезы H_0 .

Г-критерий

Г-критерий

Пусть $f_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль $F_{m,n-k}$. Если выполнено неравенство

$$\frac{\hat{Q}_A(\mathbf{X})}{\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2} \cdot \frac{n-k}{m} = \frac{(\hat{\boldsymbol{\tau}} - \boldsymbol{\tau})^T \mathbf{D}^{-1}(\hat{\boldsymbol{\tau}} - \boldsymbol{\tau})}{\|\mathbf{X} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2} \cdot \frac{n-k}{m} > f_{1-\varepsilon},$$

то гипотеза $H_0 : \mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau}$ отвергается.

Вероятность ошибки первого рода у Г-критерия равна ε .

Лемма

Г-критерий является несмещенным.

Пример: проверка однородности

Даны две независимые нормальные выборки: $\mathbf{X} = (X_1, \dots, X_n) \sim \mathcal{N}(a_1, \sigma^2)$ и $\mathbf{Y} = (Y_1, \dots, Y_m) \sim \mathcal{N}(a_2, \sigma^2)$. Постройте Г-критерий для проверки гипотезы однородности $H_0 : a_1 = a_2$.

Решение

Сведем задачу к гауссовской модели. Для этого составим единый вектор наблюдений:

$$\mathbf{W} = (X_1, \dots, X_n, Y_1, \dots, Y_m)^T = (a_1, \dots, a_1, a_2, \dots, a_2)^T + \varepsilon,$$

где $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{n+m})$. Далее, представим:

$$\begin{pmatrix} a_1 \\ \vdots \\ a_1 \\ a_2 \\ \vdots \\ a_2 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{pmatrix}}_{\mathbf{Z}} \underbrace{\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}}_{\boldsymbol{\theta}}.$$

Гауссовская линейная модель построена. Теперь перепишем гипотезу $H_0 : a_1 = a_2$ в виде $H_0 : \mathbf{A}\boldsymbol{\theta} = \mathbf{0}$, где $\mathbf{A} = \begin{pmatrix} 1 & -1 \end{pmatrix} \in \text{Mat}(1 \times 2)$. Остается вычислить участвующие величины.

1) Оценка наименьших квадратов:

$$\hat{\boldsymbol{\theta}}(\mathbf{X}) = \begin{pmatrix} \bar{\mathbf{X}} \\ \bar{\mathbf{Y}} \end{pmatrix};$$

2)

$$\mathbf{D} = \mathbf{A}(\mathbf{Z}^T \mathbf{Z})^{-1} = \frac{1}{n} + \frac{1}{m};$$

3)

$$\hat{Q}_A(\mathbf{X}) = (\bar{\mathbf{X}} - \bar{\mathbf{Y}})^2 \left(\frac{1}{n} + \frac{1}{m} \right)^{-1};$$

4)

$$\|\mathbf{W} - \mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})\|^2 = \sum_{i=1}^n (X_i - \bar{\mathbf{X}})^2 + \sum_{j=1}^m (Y_j - \bar{\mathbf{Y}})^2.$$

В итоге, F-критерий имеет вид

$$\begin{aligned} & \frac{(\bar{\mathbf{X}} - \bar{\mathbf{Y}})^2 \left(\frac{1}{n} + \frac{1}{m} \right)^{-1}}{\sum_{i=1}^n (X_i - \bar{\mathbf{X}})^2 + \sum_{j=1}^m (Y_j - \bar{\mathbf{Y}})^2} \cdot \frac{n+m-2}{1} = \\ & = \frac{(\bar{\mathbf{X}} - \bar{\mathbf{Y}})^2}{\sum_{i=1}^n (X_i - \bar{\mathbf{X}})^2 + \sum_{j=1}^m (Y_j - \bar{\mathbf{Y}})^2} \cdot \frac{nm(n+m-2)}{n+m} > f_{1-\varepsilon}, \end{aligned}$$

где $f_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль $F_{1, m+m-2}$.

Замечание

Существенным недостатком F-критерия является необходимость обращения матрицы $\mathbf{D} = \mathbf{A}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{A}^T$. Даже в самых простых постановках это бывает невозможно сделать аналитически.

Обобщенный метод наименьших квадратов

Оказывается, статистику $\hat{Q}_A(\mathbf{X})$ можно находить по-другому, используя прямую оптимизацию. Обозначим

$$s(\boldsymbol{\theta}) = \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2$$

и положим

$$\boldsymbol{\theta}^*(\mathbf{X}) = \arg \min_{\boldsymbol{\theta}: \mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau}} \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2.$$

Лемма

Выполняется равенство

$$\boldsymbol{\theta}^*(\mathbf{X}) = \hat{\boldsymbol{\theta}}(\mathbf{X}) - (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{A}^T \mathbf{D}^{-1} (\mathbf{A} \hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\tau}).$$

Доказательство.

Заметим, что

$$\mathbf{A}\boldsymbol{\theta}^*(\mathbf{X}) = \mathbf{A}\hat{\boldsymbol{\theta}}(\mathbf{X}) - \mathbf{D}\mathbf{D}^{-1}(\mathbf{A}\hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\tau}) = \boldsymbol{\tau},$$

так что $\boldsymbol{\theta}^*$ может быть искомым.

Теперь для любого $\boldsymbol{\theta} \in \mathbb{R}^k$ рассмотрим:

$$\begin{aligned} s(\boldsymbol{\theta}) &= \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2 = \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta} \pm \mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X})\|^2 = \\ &= s(\boldsymbol{\theta}^*(\mathbf{X})) + \|\mathbf{Z}\boldsymbol{\theta} - \mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X})\|^2 + 2(\mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X}) - \mathbf{Z}\boldsymbol{\theta})^T(\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X})). \end{aligned}$$

Покажем теперь, что если $\mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau}$, то последнее слагаемое равно нулю. Действительно,

$$\begin{aligned} (\mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X}) - \mathbf{Z}\boldsymbol{\theta})^T(\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X})) &= (\boldsymbol{\theta}^*(\mathbf{X}) - \boldsymbol{\theta})^T(\mathbf{Z}^T\mathbf{X} - \mathbf{Z}^T\mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X})) \quad \square \\ &\quad (\text{так как } \mathbf{Z}^T\mathbf{X} = \mathbf{Z}^T\mathbf{Z}\hat{\boldsymbol{\theta}}(\mathbf{X})) \\ &\quad \square (\boldsymbol{\theta}^*(\mathbf{X}) - \boldsymbol{\theta})^T\mathbf{Z}^T\mathbf{Z}(\hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\theta}^*(\mathbf{X})) \quad \square \\ &\quad (\text{применим к } \boldsymbol{\theta}^*(\mathbf{X}) \text{ во второй скобке формулу из условия Леммы}) \\ &\quad \square (\boldsymbol{\theta}^*(\mathbf{X}) - \boldsymbol{\theta})^T\mathbf{Z}^T\mathbf{Z}(\mathbf{Z}^T\mathbf{Z})^{-1}\mathbf{A}^T\mathbf{D}^{-1}(\mathbf{A}\hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\tau}) \quad \square \\ &\quad (\text{так как } \mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau}) \\ &\quad \square (\boldsymbol{\theta}^*(\mathbf{X}) - \boldsymbol{\theta})^T\mathbf{A}^T\mathbf{D}^{-1}\mathbf{A}(\hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\theta}^*(\mathbf{X})). \end{aligned}$$

Остается заметить, что уже первая часть равна нулю:

$$(\boldsymbol{\theta}^*(\mathbf{X}) - \boldsymbol{\theta})^T\mathbf{A}^T = (\mathbf{A}\boldsymbol{\theta}^*(\mathbf{X}) - \mathbf{A}\boldsymbol{\theta})^T = (\boldsymbol{\tau} - \boldsymbol{\tau})^T = 0.$$

Значит,

$$s(\boldsymbol{\theta}) = s(\boldsymbol{\theta}^*(\mathbf{X})) + \|\mathbf{Z}\boldsymbol{\theta} - \mathbf{Z}\boldsymbol{\theta}^*(\mathbf{X})\|^2 \geq s(\boldsymbol{\theta}^*(\mathbf{X}))$$

для всех $\boldsymbol{\theta}$ с условием $\mathbf{A}\boldsymbol{\theta} = \boldsymbol{\tau}$. То есть $\boldsymbol{\theta}^*(\mathbf{X})$ решает экстремальную задачу, на ней достигается наименьшее значение функции s на аффинном подпространстве. Причем в силу полноты ранга \mathbf{Z} равенство будет достигаться тогда и только тогда, когда $\boldsymbol{\theta} = \boldsymbol{\theta}^*(\mathbf{X})$. ■

$\boldsymbol{\theta}^*$ можно находить как решение экстремальной задачи, а также с ее помощью можно вычислять статистику \hat{Q}_A , которая нужна в критерии. Вторая лемма играет ключевую роль в обобщенном методе наименьших квадратов.

Лемма

Выполняется равенство

$$\hat{Q}_A(\mathbf{X}) = s(\boldsymbol{\theta}^*(\mathbf{X})) - s(\hat{\boldsymbol{\theta}}(\mathbf{X})).$$

Доказательство.

Из доказательства предыдущей леммы следует, что для любого $\theta \in \mathbb{R}^k$ выполняется равенство:

$$\begin{aligned} s(\theta) &= s(\theta^*(X)) + \|\mathbf{Z}\theta - \mathbf{Z}\theta^*(X)\|^2 + 2(\theta^*(X) - \theta)^T \mathbf{A}^T \mathbf{D}^{-1} (\mathbf{A}\hat{\theta}(X) - \tau) = \\ &= s(\theta^*(X)) + \|\mathbf{Z}\theta - \mathbf{Z}\theta^*(X)\|^2 - 2(\mathbf{A}\theta - \tau)^T \mathbf{D}^{-1} (\mathbf{A}\hat{\theta}(X) - \tau). \end{aligned}$$

Положим $\theta = \hat{\theta}(X)$, тогда последнее слагаемое будет в точности равно $-2 \cdot \hat{Q}_A(X)$ (так как $\hat{Q}_A(X) = (\hat{\tau} - \tau)^T \mathbf{D}^{-1} (\hat{\tau} - \tau)$). Рассмотрим второе слагаемое и попробуем показать, что оно тоже совпадает со статистикой \hat{Q}_A :

$$\|\mathbf{Z}\hat{\theta}(X) - \mathbf{Z}\theta^*(X)\|^2 = (\hat{\theta}(X) - \theta^*(X))^T \mathbf{Z}^T \mathbf{Z} (\hat{\theta}(X) - \theta^*(X)) \quad \square$$

(используем формулу для $\theta^*(X)$)

$$\begin{aligned} \square &= (\mathbf{A}\hat{\theta}(X) - \tau)^T \mathbf{D}^{-1} \mathbf{A} (\mathbf{Z}^T \mathbf{Z})^{-1} \cdot \mathbf{Z}^T \mathbf{Z} \cdot (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{A}^T \mathbf{D}^{-1} (\mathbf{A}\hat{\theta}(X) - \tau) = \\ &= (\mathbf{A}\hat{\theta}(X) - \tau)^T \mathbf{D}^{-1} \cdot \mathbf{D} \cdot \mathbf{D}^{-1} (\mathbf{A}\hat{\theta}(X) - \tau) = (\mathbf{A}\hat{\theta}(X) - \tau)^T \mathbf{D}^{-1} (\mathbf{A}\hat{\theta}(X) - \tau) = \hat{Q}_A(X). \end{aligned}$$

В итоге, $s(\hat{\theta}(X)) = s(\theta^*(X)) - \hat{Q}_A(X)$. ■

Вывод

Для вычисления статистики \hat{Q}_A достаточно найти $s(\hat{\theta}(X))$ и $s(\theta^*(X))$, которые могут быть найдены путем решения экстремальных задач методами оптимизации. Кроме того, итоговая статистика для применения F-критерия будет иметь вид

$$\frac{s(\theta^*(X)) - s(\hat{\theta}(X))}{s(\hat{\theta}(X))} \cdot \frac{n - k}{m}.$$

Пример: проверка однородности

Даны k независимых нормальных выборок:

$$\mathbf{X}^{(1)} = (X_1^{(1)}, \dots, X_{n_1}^{(1)}) \sim \mathcal{N}(a_1, \sigma^2), \dots, \mathbf{X}^{(k)} = (X_1^{(k)}, \dots, X_{n_k}^{(k)}) \sim \mathcal{N}(a_k, \sigma^2).$$

Постройте F-критерий для проверки гипотезы однородности $H_0 : a_1 = a_2 = \dots = a_k$.

Решение

Сведем задачу к гауссовской модели. Для этого составим единый вектор наблюдений:

$$\mathbf{W} = (X_1^{(1)}, \dots, X_{n_1}^{(1)}, X_1^{(2)}, \dots, X_{n_2}^{(2)}, \dots, X_1^{(k)}, \dots, X_{n_k}^{(k)})^T = (a_1, \dots, a_1, a_2, \dots, a_k)^T + \varepsilon,$$

где $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{n_1 + \dots + n_k})$. Представляя вектор среднего в виде $\mathbf{Z}\theta$, где $\theta = (a_1, a_2, \dots, a_k)^T$, получаем, что

$$s(\theta) = \|\mathbf{W} - \mathbf{Z}\theta\|^2 = \sum_{j=1}^k \sum_{i=1}^{n_j} (X_i^{(j)} - a_j)^2.$$

Далее,

$$s(\boldsymbol{\theta}^*(\mathbf{X})) = \min_{\boldsymbol{\theta} \in \mathbb{R}^k} \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2 = \sum_{j=1}^k \sum_{i=1}^{n_j} (X_i^{(j)} - \overline{\mathbf{X}}^{(j)})^2,$$

где $\overline{\mathbf{X}}^{(j)} = \frac{1}{n_j} \sum_{i=1}^{n_j} X_i^{(j)}$.

В свою очередь, минимизируя $s(\boldsymbol{\theta})$ по подпространству $a_1 = a_2 = \dots = a_k$, получаем:

$$\begin{aligned} s(\boldsymbol{\theta}^*(\mathbf{X})) &= \min_{\boldsymbol{\theta} \in \mathbb{R}^k : a_1 = a_2 = \dots = a_k} \|\mathbf{X} - \mathbf{Z}\boldsymbol{\theta}\|^2 = \\ &= \min_{a \in \mathbb{R}} \sum_{j=1}^k \sum_{i=1}^{n_j} (X_i^{(j)} - a)^2 = \sum_{j=1}^k \sum_{i=1}^{n_j} (X_i^{(j)} - \overline{\mathbf{X}})^2, \end{aligned}$$

где $\overline{\mathbf{X}} = \frac{1}{n_1 + \dots + n_k} \sum_{j=1}^k \sum_{i=1}^{n_j} X_i^{(j)}$.

Тем самым, мы нашли $s(\boldsymbol{\theta}^*(\mathbf{X}))$ и, соответственно, $s(\hat{\boldsymbol{\theta}}(\mathbf{X}))$. Осталось применить формулу из вывода, чтобы получить статистику для F-критерия и сравнить ее с квантилью.

На этом мы завершаем разбор проверки линейных гипотез в линейной гауссовской модели.

17.

Лекция 17

Критерии согласия

Рассмотрим так называемые критерии согласия.

Определение 17.1. Пусть \mathbf{X} - наблюдение с неизвестным распределением P . Далее, введем гипотезы $H_0 : P = P_0$ и $H_1 : P \neq P_0$. Критерий для проверки гипотезы H_0 против полной альтернативы H_1 называется критерием согласия.

Что мы хотели бы от такого критерия? В идеале мы желаем получить что-то наподобие равномерно наиболее мощного критерия. Но эта задача почти всегда безнадежна. Даже в случае хороших распределений, для которых выполнено свойство монотонного отношения правдоподобия, 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.

Упражнение

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из $U(0, \theta)$. Постройте р.н.м.к. уровня значимости ε для проверки гипотезы $H_0 : \theta = \theta_0$ против альтернативы $H_1 : \theta \neq \theta_0$.

Асимптотические критерии согласия

Совершенно естественно на этом пути возникает задача поиска асимптотических критериев. Что мы понимаем под асимптотическим критерием? Как и в случае с оценками, если мы говорим об асимптотике, то она берется по размеру выборки. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка растущего размера (для каждого n мы добавляем еще один элемент выборки). 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 $H_1^\theta : \theta \neq \theta_0$ такую, что

- Предел вероятности ошибки первого рода стремится к $\varepsilon \in (0, 1)$:

$$P_{\theta_0}(\mathbf{X} \in S_n) \rightarrow \varepsilon.$$

- Последовательность критериев будет состоятельна: для любого $\theta \neq \theta_0$

$$P_\theta(\mathbf{X} \in S_n) \rightarrow 1.$$

Несмотря на то, что условия кажутся сложными для выполнения, оказывается, что решения существуют. Рассмотрим первый из таких критериев - критерий хи-квадрат Пирсона для случая дискретных распределений.

Критерий хи-квадрат

Пусть есть выборка $\mathbf{X} = (X_1, \dots, X_n)$ из распределения P с конечным числом значений $\{a_1, \dots, a_m\}$:

$$p_j = P(X_1 = a_j) > 0, \quad j = 1, \dots, m, \quad \sum_{j=1}^m p_j = 1.$$

Введем вектор вероятностей $\mathbf{p} = (p_1, \dots, p_m)$. Гипотеза $H_0 : P = P_0$ тогда сведется к равенству векторов вероятностей: $H_0 : \mathbf{p} = \mathbf{p}^0$. Альтернатива примет вид $H_1 : \mathbf{p} \neq \mathbf{p}^0$.

Для каждого $j = 1, \dots, m$ введем

$$\mu_j = \sum_{i=1}^n \mathbf{I}\{X_i = a_j\},$$

- это число элементов выборки, равных a_j .

Определение 17.2. Статистикой хи-квадрат для проверки гипотезы $H_0^{\mathbf{p}} = \mathbf{p}^0$ называется

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{j=1}^m \frac{(\mu_j - np_j^0)^2}{np_j^0}.$$

Сам же критерий устроен следующим образом.

Критерий хи-квадрат

Пусть $u_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль χ_{m-1}^2 . Если выполнено равенство

$$\hat{\chi}_n^2(\mathbf{X}) > u_{1-\varepsilon},$$

то гипотеза $H_0 : \mathbf{p} = \mathbf{p}^0$ отвергается.

Изучим свойство этого критерия, а именно - найдем предельные вероятности ошибок первого и второго рода. Начнем со второго.

Утверждение

Критерий хи-квадрат состоятелен.

Доказательство.

Для начала немного перепишем статистику хи-квадрат:

$$\hat{\chi}_n^2(\mathbf{X}) = n \sum_{j=1}^m \frac{1}{p_j^0} \left(\frac{\mu_j}{n} - p_j^0 \right)^2.$$

Далее, если $\mathbf{p} \neq \mathbf{p}^0$, то найдется индекс $j \in \{1, 2, \dots, m\}$ такой, что $p_j \neq p_j^0$. По усиленному закону больших чисел $\frac{\mu_j}{n} \rightarrow p_j$ почти наверное. Отсюда следует, что $\hat{\chi}_n^2(\mathbf{X})$ линейно стремится к бесконечности, так как

$$\frac{\hat{\chi}_n^2(\mathbf{X})}{n} \xrightarrow{\text{п.н.}} \sum_{j=1}^m \frac{(p_j - p_j^0)^2}{p_j^0} > 0.$$

Тем самым, критерий состоятелен: $P(\hat{\chi}_n^2(\mathbf{X}) > u_{1-\varepsilon}) \rightarrow 1$. ■

Разберемся теперь, что происходит с вероятностью ошибки первого рода, то есть что происходит со статистикой, если гипотеза правильная. Название критерия и его построение вытекает из следующей теоремы

Теорема 17.1. (Пирсон)

Если выполнена гипотеза H_0 , то имеет место следующая сходимост по распределению:

$$\hat{\chi}_n^2(\mathbf{X}) \xrightarrow[n \rightarrow \infty]{d} \chi_{m-1}^2.$$

Доказательство.

Доказательство состоит в работе с гауссовскими случайными векторами, а распределение хи-квадрат возникает при попытке перейти от многомерного варианта к одномерному. Распределение χ_{m-1}^2 получается как сумма квадратов стандартных нормальных случайных величин. Поймем теперь, где они возникают.

Пусть $\mathbf{p} = \mathbf{p}^0$. Для всех $i = 1, \dots, n$ введем случайный вектор

$\mathbf{Y}_i = (\mathbf{I}\{X_i = a_1\}, \dots, \mathbf{I}\{X_i = a_m\})^T$. Понятно, что такие векторы независимы и одинаково распределены, причем $E\mathbf{Y}_i = \mathbf{p}^0$. Далее, заметим, что

$$\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} = \frac{1}{n} \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}.$$

Следовательно, по многомерной центральной предельной теореме

$$\sqrt{n} \left(\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} - \mathbf{p}^0 \right) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(\mathbf{0}, \Sigma)$$

где $\Sigma = D\mathbf{Y}_1$. - матрица ковариаций случайного вектора \mathbf{Y}_1 . Заметим, что

$$\text{cov}(\mathbf{I}\{X_1 = a_i\}, \mathbf{I}\{X_1 = a_j\}) = P(X_1 = a_i, X_1 = a_j) - p_i^0 p_j^0 = \begin{cases} p_i^0 - p_i^0 p_j^0, & i = j; \\ -p_i^0 p_j^0, & i \neq j. \end{cases}$$

Неприятность состоит в том, что это вырожденная матрица. Запишем матрицу ковариаций следующим образом: $\Sigma = \mathbf{B} - \mathbf{p}^0(\mathbf{p}^0)^T$, где $\mathbf{B} = \text{diag}(\mathbf{p}^0)$. Теперь наша

цель - повернуть вектор так, чтобы матрица ковариаций стала похожа на матрицу стандартного нормального распределения. Положим

$$\xi'_n = \sqrt{n} \left(\frac{\mathbf{Y}_1 + \dots + \mathbf{Y}_n}{n} - \mathbf{p}^0 \right).$$

Тогда по теореме о наследовании сходимости

$$(\sqrt{\mathbf{B}})^{-1} \xi'_n \xrightarrow[n \rightarrow \infty]{d} (\sqrt{\mathbf{B}})^{-1} \mathcal{N}(\mathbf{0}, \Sigma) = \mathcal{N}(\mathbf{0}, (\sqrt{\mathbf{B}})^{-1} \Sigma (\sqrt{\mathbf{B}})^{-1}).$$

Новая матрица ковариаций равна $\mathbf{I}_m - \mathbf{z}\mathbf{z}^T$, где

$$\mathbf{z} = (\sqrt{\mathbf{B}})^{-1} \mathbf{p}^0 = \left(\sqrt{p_1^0}, \dots, \sqrt{p_m^0} \right)^T.$$

Заметим, что вектор \mathbf{z} единичный, значит, он может быть базисным. Рассмотрим ортогональную матрицу \mathbf{C} , у которой первая строка равна \mathbf{z} , а остальные равны чему угодно. Тогда по теореме о наследовании сходимости

$$\xi''_n = \mathbf{C}(\sqrt{\mathbf{B}})^{-1} \xi'_n \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(\mathbf{0}, \mathbf{C}(\mathbf{I}_m - \mathbf{z}\mathbf{z}^T)\mathbf{C}^T).$$

Однако $\mathbf{C}\mathbf{z} = (1, 0, \dots, 0)^T$, а $\mathbf{C}\mathbf{C}^T = \mathbf{I}_m$. Поэтому

$$\mathbf{C}(\mathbf{I}_m - \mathbf{z}\mathbf{z}^T)\mathbf{C}^T = \text{diag}(0, 1, \dots, 1) = \mathbf{I}'_m.$$

Снова по теореме о наследовании сходимости

$$\|\xi''_n\|^2 \xrightarrow[n \rightarrow \infty]{d} \|\mathcal{N}(\mathbf{0}, \mathbf{I}'_m)\|^2 \sim \chi_{m-1}^2.$$

Но, как известно, умножение вектора на ортогональную матрицу не меняет его норму. Тогда

$$\|(\sqrt{\mathbf{B}})^{-1} \xi'_n\|^2 \xrightarrow[n \rightarrow \infty]{d} \chi_{m-1}^2.$$

Осталось заметить, что

$$(\sqrt{\mathbf{B}})^{-1} \xi'_n = (\sqrt{\mathbf{B}})^{-1} \sqrt{n} \begin{pmatrix} \frac{\mu_1 - np_1^0}{n} \\ \dots \\ \frac{\mu_m - np_m^0}{n} \end{pmatrix} = \begin{pmatrix} \frac{\mu_1 - np_1^0}{\sqrt{np_1^0}} \\ \dots \\ \frac{\mu_m - np_m^0}{\sqrt{np_m^0}} \end{pmatrix}.$$

Тем самым,

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{j=1}^m \frac{(\mu_j - np_j^0)^2}{np_j^0} = \|(\sqrt{\mathbf{B}})^{-1} \xi'_n\|^2 \xrightarrow[n \rightarrow \infty]{d} \chi_{m-1}^2.$$

■

Следствие

Для критерия хи-квадрат вероятность ошибки первого рода стремится к ε с ростом n .

Пожертвовав точностью вероятности ошибки первого рода, мы получили критерий с очень сильными асимптотическими свойствами.

Практическая применимость

Насколько этот метод применим? На практике считается, что должно быть выполнено следующее ограничение: $np_j^0 \geq 5$ для всех $j = 1, \dots, m$.

Параметрический критерий хи-квадрат

Предположим, что мы хотим проверить более сложную гипотезу. Оказывается, что критерий хи-квадрат можно естественным образом обобщить. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из неизвестного распределения со значениями $\{a_1, \dots, a_m\}$:

$$p_i(\theta) = P_\theta(X_1 = a_i), \theta \in \Theta.$$

В данном случае гипотезу формулируется в виде принадлежности к параметрическому семейству: $H_0 : \mathbf{p} \in \{\mathbf{p}_\theta, \theta \in \Theta\}$ против $H_1 : \mathbf{p} \notin \{\mathbf{p}_\theta, \theta \in \Theta\}$. Снова хотелось бы составить статистику хи-квадрат, но вектора \mathbf{p}^0 больше нет. Можно ли его чем-то заменить?

Естественно заменить \mathbf{p}^0 на некоторую оценку:

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{i=1}^m \frac{(\mu_i - n\hat{p}_i(\mathbf{X}))^2}{n\hat{p}_i(\mathbf{X})}. \quad (1)$$

Какую оценку взять? Самым разумным выбором является оценка максимального правдоподобия: $\hat{p}_i(\mathbf{X}) = p_i(\hat{\theta}(\mathbf{X}))$, где $\hat{\theta}(\mathbf{X})$ - оценка максимального правдоподобия для θ . Оказывается, что в таком случае можно доказать результат, похожий на теорему Пирсона.

Условия регулярности

Напомним, что ОМП вычисляется как решение задачи

$$\hat{\theta}_n(\mathbf{X}) = \arg \max_{\theta \in \Theta} \prod_{i=1}^m (p_i(\theta))^{\mu_i}.$$

Как всегда, если мы имеем дело с ОМП, то необходимы условия регулярности. Пусть выполнены следующие условия регулярности.

- 1) $\Theta \subseteq \mathbb{R}^s$, $s < m - 1$ - открытое множество.
- 2) Для любого значения параметра все вероятности отделены от нуля: $p_i(\theta) \geq c^2 > 0$ для всех i и θ .
- 3) 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 всем Θ .
- 4) Матрица $\mathbf{D}(\theta) = \left(\frac{\partial p_i(\theta)}{\partial \theta_j}, i = 1, \dots, m, j = 1, \dots, s \right)$ имеет ранг s для любого $\theta \in \Theta$.

Теперь можно сформулировать параметрический аналог теоремы Пирсона.

Теорема 17.2. (Фишер)

Пусть выполнены условия регулярности. Введем следующую систему уравнений:

$$\sum_{i=1}^m \frac{\mu_i}{p_i(\theta)} \frac{\partial p_i(\theta)}{\partial \theta_j} = 0, \quad j = 1, \dots, s, \quad \mu_j = \sum_{i=1}^n \mathbf{I}\{X_i = a_j\}.$$

Если верна гипотезы H_0 , то с вероятностью, стремящейся к 1, данная система имеет единственное решение $\hat{\theta}_n(\mathbf{X})$ такое, что $\hat{\theta}_n(\mathbf{X})$ сходится по вероятности к истинному значению параметра θ_0 и

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{i=1}^m \frac{(\mu_i - np_i(\hat{\theta}_n(\mathbf{X})))^2}{np_i(\hat{\theta}_n(\mathbf{X}))} \xrightarrow[n \rightarrow \infty]{d_{\theta_0}} \chi_{m-1-s}^2.$$

Опустим доказательство этой теоремы и посмотрим на ее применение.

Параметрический критерий хи-квадрат

Пусть $u_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль χ_{m-1-s}^2 . Если выполнено неравенство

$$\hat{\chi}_n^2(\mathbf{X}) > u_{1-\varepsilon},$$

то гипотеза $H_0 : \mathbf{p} \in \{\mathbf{p}_\theta, \theta \in \Theta\}$ отвергается.

Рассмотрим применение теоремы и параметрического критерия хи-квадрат к построению критерия независимости хи-квадрат.

Критерий независимости хи-квадрат

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из векторного распределения, $X_k = (Y_k, Z_k)$, где случайные величины Y_k принимают значения A_1, \dots, A_s , а Z_k - значения B_1, \dots, B_l ("признаки"). Это так называемая двухфакторная модель.

Хотим проверить гипотезу независимости признаков: H_0 : признаки независимы.

Обозначим

$$p_i = P(Y_1 = A_i), \quad i = 1, \dots, s; \quad q_j = P(Z_1 = B_j), \quad j = 1, \dots, l;$$

$$p_{ij} = P(X_1 = (A_i, B_j)), \quad i = 1, \dots, s; \quad j = 1, \dots, l.$$

Сведем задачу к параметрической и построим параметрический критерий для проверки гипотезы H_0 против полной альтернативы H_1 : признаки зависимы.

Введем параметр $\theta = (p_1, \dots, p_{s-1}, q_1, \dots, q_{l-1})$ с условиями

$$0 < p_i, q_i, \quad i = 1, \dots, s, \quad j = 1, \dots, l; \quad \sum_{i=1}^{s-1} p_i < 1, \quad \sum_{j=1}^{l-1} q_j < 1.$$

Размерность параметра равна $s + l - 2$. Положим $p_s = 1 - \sum_{i=1}^{s-1} p_i$, $q_l = 1 - \sum_{j=1}^{l-1} q_j$. В параметрическом семействе, отвечающем гипотезе H_0 , выполнено

$$p_{ij} = p_i \cdot q_j, \quad i = 1, \dots, s; \quad j = 1, \dots, l.$$

Для построения статистики $\hat{\chi}_n^2(\mathbf{X})$ нам необходимо определить остальные величины.

1) Обозначим

$$\mu_{ij} = \sum_{k=1}^n \mathbf{I}\{X_k = (A_i, B_j)\}, \quad \mu_i = \sum_{k=1}^n \mathbf{I}\{Y_k = A_i\}, \quad \nu_j = \sum_{k=1}^n \mathbf{I}\{Z_k = B_j\}.$$

2) Найдем оценки максимального правдоподобия для параметра $\theta = (p_1, \dots, p_{s-1}, q_1, \dots, q_{l-1})$.

Заметим, что функция правдоподобия выборки в случае вероятности гипотезы равна

$$f_{\theta}(\mathbf{X}) = \prod_{i=1}^s \prod_{j=1}^l (p_i \cdot q_j)^{\mu_{ij}} = \prod_{i=1}^s p_i^{\mu_i} \cdot \prod_{j=1}^l q_j^{\nu_j}.$$

Тем самым, оптимизация дальше идет отдельно по (p_1, \dots, p_{s-1}) и (q_1, \dots, q_{l-1}) .

Утверждение

ОМП для (p_1, \dots, p_{s-1}) равны

$$\hat{p}_i(\mathbf{X}) = \frac{\mu_i}{n}, \quad i = 1, \dots, s-1.$$

Доказательство.

Возьмем частную производную по p_i у логарифма $f_{\theta}(\mathbf{X})$:

$$\frac{\partial}{\partial p_i} \ln f_{\theta}(\mathbf{X}) = \frac{\mu_i}{p_i} - \frac{\mu_s}{1 - \sum_{k=1}^{s-1} p_k}.$$

Равенство нулю всех производных означает, что все величины $\frac{p_i}{\mu_i}$ равны некоторому одному и тому же α . Из условия $\sum_{k=1}^s p_k = 1$ получаем, что

$$1 = \sum_{k=1}^s p_k = \alpha \sum_{k=1}^s \mu_k = \alpha \cdot n.$$

Значит, $\hat{p}_i(\mathbf{X}) = \frac{\mu_i}{n}$, $i = 1, \dots, s$. ■

В итоге обобщенная статистика хи-квадрат равна

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{i=1}^s \sum_{j=1}^l \frac{(\mu_{ij} - n\hat{p}_i(\mathbf{X})\hat{q}_j(\mathbf{X}))^2}{n\hat{p}_i(\mathbf{X})\hat{q}_j(\mathbf{X})} = \sum_{i=1}^s \sum_{j=1}^l \frac{(n\mu_{ij} - \mu_i \cdot \nu_j)^2}{n \cdot \mu_i \cdot \nu_j}.$$

3) Число степеней свободы равно $sl - 1 - (s + l - 2) = (s - 1)(l - 1)$.

Критерий независимости хи-квадрат

Пусть $u_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль $\chi_{(s-1)(l-1)}^2$. Если выполнено неравенство

$$\hat{\chi}_n^2(\mathbf{X}) > u_{1-\varepsilon},$$

то гипотеза H_0 : отвергается.

По теореме Фишера данный критерий имеет вероятность ошибки первого рода в пределе равную ε .

Упражнение

Докажите, что критерий независимости состоятелен.

Критерий однородности хи-квадрат

Еще одним обобщением стандартного критерия хи-квадрат является критерий однородности хи-квадрат. В отличие от критерия независимости, критерий однородности не выводится из параметрического случая.

Задача состоит в проверке однородности выборок. Пусть

$$\mathbf{X}^{(1)} = (X_1^1, \dots, X_{n_1}^1), \dots, \mathbf{X}^{(k)} = (X_k^1, \dots, X_k^{n_k})$$

- независимые выборки, принимающие одни и те же значения $\{a_1, \dots, a_m\}$. Обозначим

$$\mathbf{p}^{(j)} = (p_1^{(j)}, \dots, p_m^{(j)}), \quad p_i^{(j)} = P(X_1^{(j)} = a_i), \quad j = 1, \dots, k.$$

Мы хотим проверить гипотезу однородности

$$H_0 : \mathbf{p}^{(1)} = \dots = \mathbf{p}^{(k)}.$$

Для этого также составляется обобщенная статистика хи-квадрат. Для этого обозначим:

$$\mu_{ij} = \sum_{s=1}^{n_j} \mathbf{I}\{X_s^{(j)} = a_i\}, \quad \mu_i = \sum_{j=1}^k \mu_{ij}, \quad n = n_1 + \dots + n_k.$$

Тогда обобщенная статистика хи-квадрат равна

$$\hat{\chi}_n^2(\mathbf{X}) = \sum_{i=1}^m \sum_{j=1}^k \frac{(\mu_{ij} - n_j \cdot \frac{\mu_i}{n})^2}{n_j \cdot \frac{\mu_i}{n}}.$$

Критерий однородности хи-квадрат

Пусть $u_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль $\chi_{(k-1)(m-1)}^2$. Если выполнено неравенство

$$\hat{\chi}_n^2(\mathbf{X}) > u_{1-\varepsilon},$$

то гипотеза H_0 : отвергается.

Упражнение

Докажите, что критерий однородности состоятелен.

Критерии для непрерывных распределений

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из непрерывного распределения с функцией распределения $F(x)$. Мы хотели бы проверить простую гипотезу $H_0 : F = F_0$, где F_0 - заданная непрерывная функция распределения на \mathbb{R} , против полной альтернативы $H_1 : F \neq F_0$.

Здравая идея

Надо дискретизировать выборку $\mathbf{X} = (X_1, \dots, X_n)$. Разобьем прямую на полуинтервалы $\Delta_1, \dots, \Delta_m$, $\Delta_j = (a_j, b_j]$ и обозначим

$$p_i^0 = F_0(b_j) - F_0(a_j).$$

Введем выборку $\mathbf{Y} = (Y_1, \dots, Y_n)$ со значениями $\{1, \dots, m\}$ по правилу

$$Y_j = s \iff X_j \in \Delta_s.$$

Далее, проверяем критерием хи-квадрат гипотезу о том, что выборка \mathbf{Y} имеет распределение $\mathbf{p}^0 = (p_1^0, \dots, p_m^0)$.

Однако, этот метод "склеивает" те функции распределения F , которые имеют один и тот же вектор \mathbf{p}^0 для одного и того же набора полуинтервалов Δ_j . Получается, что критерий не состоятелен, но он различит содержательно те распределения, которые имеют заведомо неодинаковые векторы \mathbf{p}^0 для одних и тех же полуинтервалов.

Теорема Колмогорова

Построим по выборке эмпирическую функцию распределения:

$$\hat{F}_n(x) = \frac{1}{n} \sum_{j=1}^n \mathbf{I}\{X_j \leq x\}.$$

Следующая теорема играет важнейшую роль в статистике для обоснования применения критерия Колмогорова.

Теорема 17.3. (Колмогоров)

Если $F(x)$ непрерывна, то распределение случайной величины

$$D_n = \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)|$$

не зависит от $F(x)$. Более того,

$$\sqrt{n} \cdot D_n \xrightarrow{d} K,$$

где K - распределение Колмогорова с ф.р. $K(y) = \sum_{k \in \mathbb{Z}} (-1)^k e^{-2k^2 y^2}$, $y > 0$.

Доказывать теорему Колмогорова мы будем на следующей лекции. Пока что построим критерий Колмогорова и обсудим его свойства.

Критерий Колмогорова

Итак, снова пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из непрерывного распределения с функцией распределения $F(x)$. Мы хотели бы проверить простую гипотезу $H_0 : F = F_0$, где F_0 - заданная непрерывная функция распределения на \mathbb{R} , против полной альтернативы $H_1 : F \neq F_0$. Введем статистику

$$\hat{D}_n = \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F_0(x)|.$$

Критерий Колмогорова

Пусть $k_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль распределения Колмогорова K . Если выполнено неравенство

$$\sqrt{n} \cdot \hat{D}_n > k_{1-\varepsilon},$$

то гипотеза $H_0 : F = F_0$ отвергается.

Следствие

Для критерия Колмогорова вероятность ошибки первого рода стремится к ε с ростом n .

Утверждение

Критерий Колмогорова состоятелен

Доказательство.

Пусть настоящая функция распределения $F(x)$ не равна $F_0(x)$. По теореме Гливленко-Кантелли мы знаем, что с вероятностью 1

$$\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \rightarrow 0.$$

Следовательно, обозначив $\sup_{x \in \mathbb{R}} |F(x) - F_0(x)| = \varepsilon_0$, мы получаем, что

$$P\left(\hat{D}_n > \frac{\varepsilon_0}{2}\right) \rightarrow 1 \quad \text{при } n \rightarrow +\infty.$$

Но тогда

$$\sqrt{n} \cdot \hat{D}_n \xrightarrow{P} +\infty.$$

В итоге, получаем, что

$$P\left(\sqrt{n} \cdot \hat{D}_n > k_{1-\varepsilon}\right) \rightarrow 1 \quad \text{при } n \rightarrow +\infty.$$

■

Критерий Колмогорова-Смирнова

Для проверки однородности аналогом теоремы Колмогорова является теорема Смирнова.

Теорема 17.4. (Смирнов)

Если гипотеза однородности верна и функция распределения $F(x)$ выборок непрерывна, то распределение статистики $\hat{D}_{n,m} = \hat{D}_{n,m}(\mathbf{X}, \mathbf{Y})$ зависит только от n и m и, кроме того,

$$\sqrt{\frac{nm}{n+m}} \cdot \hat{D}_{n,m} \xrightarrow{d} K, \quad n, m \rightarrow +\infty.$$

Опираясь на эту теорему можно сформулировать критерий однородности Колмогорова-Смирнова.

Критерий однородности Колмогорова-Смирнова

Пусть $k_{1-\varepsilon}$ - это $(1 - \varepsilon)$ -квантиль распределения Колмогорова K . Если выполнено неравенство

$$\sqrt{\frac{nm}{n+m}} \cdot \hat{D}_n > k_{1-\varepsilon},$$

то гипотеза однородности $H_0 : X_i \stackrel{d}{=} Y_j$ отвергается.

Следствие

Для критерия Колмогорова-Смирнова вероятность ошибки первого рода стремится к ε с ростом n .

Замечание

Заметим, что квантили распределения Колмогорова весьма невелики, например $K(1.8) = 0.996932 \dots$

Упражнение

Докажите, что критерий однородности Колмогорова-Смирнова состоятелен.

На следующей лекции мы будем разбирать доказательство теоремы Колмогорова, которое опирается на сходимости по распределению случайных процессов. Нам нужно будет вспомнить некоторый материал из теории случайных процессов.

18.

Лекция 18

Напоминание из курса теории вероятностей

Слабая сходимость вероятностных пар

Пусть (S, ρ) - метрическое пространство.

Определение 18.1. Борелевской сигма-алгеброй, $\mathcal{B}(S)$, на (S, ρ) называется минимальная σ -алгебра, содержащая все открытые множества в S .

Определение 18.2. Пусть задано метрическое пространство (S, ρ) и последовательность $\{Q_n, n \in \mathbb{N}\}$ вероятностных мер на $(S, \mathcal{B}(S))$. Будем говорить, что Q_n слабо сходятся к вероятностной мере Q на $(S, \mathcal{B}(S))$, если для любой ограниченной непрерывной функции $f: S \rightarrow \mathbb{R}$ выполнено

$$\lim_{n \rightarrow \infty} \int_S f(x) Q_n(dx) = \int_S f(x) Q(dx).$$

Обозначение: $Q_n \xrightarrow{\omega} Q$.

Теорема Александрова

Утверждение

Если пространство (S, ρ) сепарабельно, то $\mathcal{B}(S)$ является минимальной σ -алгеброй, содержащей все открытые шары.

Теорема 18.1. (А.Д. Александров)

Пусть $\{Q_n, n \in \mathbb{N}\}$ и Q - вероятностные меры на метрическом пространстве (S, ρ) . Тогда следующие утверждения эквивалентны:

- 1) $Q_n \xrightarrow{\omega} Q$,
- 2) $\overline{\lim}_{n \rightarrow \infty} Q_n(F) \leq Q(F)$ для любого замкнутого множества $F \subset S$,
- 3) $\underline{\lim}_{n \rightarrow \infty} Q_n(G) \geq Q(G)$ для любого открытого множества $G \subset S$,
- 4) Для любого борелевского множества $B \in \mathcal{B}(S)$ такого, что $Q(\partial B) = 0$, выполнено $Q_n(B) \rightarrow Q(B)$ при $n \rightarrow \infty$.

Сходимость случайных процессов

Будем рассматривать время $T = [0, 1]$.

Пусть $X^n = (X_t^{(n)}, t \in [0, 1], n \in \mathbb{N})$ - последовательность действительных случайных процессов на $[0, 1]$. Как можно определить сходимость по распределению такой последовательности?

Идея 1: как слабую сходимость их распределений!

Вопрос: а что такое распределение случайного процесса?

Для этого мы напомним основные определения случайных процессов.

Случайные процессы

Определение 18.3. Пусть (Ω, \mathcal{F}, P) - вероятностное пространство. Отображение $\xi : \Omega \rightarrow \mathbb{R}$ называется случайной величиной, если оно измеримо, то есть

$$\forall B \in \mathcal{B}(\mathbb{R}) \quad \xi^{-1}(B) = \{\omega : \xi(\omega) \in B\} \in \mathcal{F}.$$

Определение 18.4. Пусть T - некоторое множество. Тогда набор $X = (X_t, t \in T)$ случайных величин $X_t(\omega)$, заданных на одном и том же вероятностном пространстве (Ω, \mathcal{F}, P) для $\forall t \in T$, называется случайным процессом с множеством времени T .

- 1) Если $T = [a, b], (a, b), [a, +\infty], \mathbb{R}$, то процесс X называется процессом с непрерывным временем.
- 2) Если $T \subset \mathbb{Z}$, то процесс X называется процессом с дискретным временем.
- 3) Если $T \subset \mathbb{R}^d, d > 1$, то процесс X называется случайным полем.

Пространство траекторий

Определение 18.5. При фиксированном $\omega = \omega_0$ функция

$$\tilde{X}_{\omega_0}(t) = X_t(\omega) \Big|_{\omega=\omega_0}$$

на T называется траекторией или реализацией случайного процесса $X = (X_t, t \in T)$.

Определение 18.6. Множество $S = \prod_{t \in T} \mathbb{R} = \mathbb{R}^T$ называется пространством траекторий случайного процесса X (здесь $\prod_{t \in T} \mathbb{R}$ -декартово произведение), то есть формально

$$S = \{y = (y(t), t \in T) : \forall t \in T \quad y(t) \in \mathbb{R}\}.$$

Пространство траекторий - это множество значений случайного процесса.

Цилиндрическая сигма-алгебра

Определение 18.7. Для $\forall t \in T$ и $B_t \in \mathcal{B}(\mathbb{R})$ введем элементарный цилиндр с основанием B_t :

$$C(t, B_t) = \{y \in S : y(t) \in B_t\}.$$

Образно говоря, $C(t, B_t)$ состоит из тех функций $y(t)$, которые в точке t проходят через множество B_t (см. рис. 6).

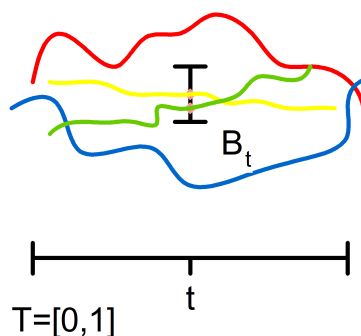


Рис. 6

Определение 18.8. Минимальная σ -алгебра \mathcal{B}_T , содержащая все эти элементарные цилиндры, называется цилиндрической σ -алгеброй на S . Формально,

$$\mathcal{B}_T = \sigma\{C(t, B_t) : t \in T, B_t \in \mathcal{B}(\mathbb{R})\}.$$

Для \mathcal{B}_T используется также обозначение $\bigotimes_{t \in T} \mathcal{B}(\mathbb{R})$.

Измеримые отображения

Стало быть, (S, \mathcal{B}_T) - измеримое пространство. Значит, на $X = (X_t, t \in T)$ можно смотреть как на один случайный элемент со значениями в (S, \mathcal{B}_T) : каждому $\omega \in \Omega$ сопоставляется целая траектория $\tilde{X}_\omega(\cdot)$. Разберемся с вопросом измеримости такого отображения.

Лемма (об измеримости случайного процесса)

$X = (X_t, t \in T)$ является случайным процессом на (Ω, \mathcal{F}, P) , то есть семейство $X = (X_t, t \in T)$ есть семейство случайных величин \Leftrightarrow отображение $X : \Omega \rightarrow S$

является измеримым относительно цилиндрической σ -алгебры.

Для доказательства нам понадобится.

Лемма (достаточное условие измеримости отображения)

Пусть (Ω, \mathcal{F}, P) - вероятностное пространство, (E, \mathcal{E}) - измеримое пространство, $\mathcal{M} \subset \mathcal{E}$ - подсистема такая, что $\sigma(\mathcal{M}) = \mathcal{E}$. Пусть $X : \Omega \rightarrow E$. Если $\forall B \in \mathcal{M} X^{-1}(B) \in \mathcal{F}$, то отображение X измеримо, то есть

$$\forall B \in \mathcal{E} X^{-1}(B) \in \mathcal{F}.$$

Доказательство.(леммы об измеримости случайного процесса)

(\Leftarrow) Надо доказать, что $\forall t X_t$ - случайная величина. Для $\forall B_t \in \mathcal{B}(\mathbb{R})$:

$$X_t^{-1}(B_t) = \{\omega : X_t(\omega) \in B_t\} = \{\omega : X(\omega) \in C(t, B_t)\} = X^{-1}(C(t, B_t)) \in \mathcal{F},$$

так как X измерим.

(\Rightarrow) Воспользуемся достаточным условием измеримости, взяв в качестве \mathcal{M} множество всех элементарных цилиндров.

Тогда $\sigma(\mathcal{M}) = \mathcal{B}_T$ - цилиндрическая σ -алгебра,

$$X^{-1}(C(t, B_t)) = \{\omega : X(\omega) \in C(t, B_t)\} = \{\omega : X_t(\omega) \in B_t\} = X_t^{-1}(B_t) \in \mathcal{F},$$

так как X_t - случайная величин. ■

Распределение случайного процесса

Замечание

Лемма устанавливает эквивалентное определение случайного процесса как единого случайного элемента со значениями в пространстве траекторий, измеримого относительно цилиндрической σ -алгебры.

Определение 18.9. Распределение P_X случайного процесса $X = (X_t, t \in T)$ называется вероятностная мера на (S, \mathcal{B}_T) такая, что

$$\forall C \in \mathcal{B}_T P_X(C) = P(X \in C).$$

Это определение осязаемо только в том случае, когда "время" конечно и мы фактически имеем дело со случайным вектором. Для счетного (или тем более континуального) "времени" это определение весьма непросто для понимания в связи с отсутствием функции распределения.

Сходимость случайных процессов

Вернемся к сходимостям. Будем рассматривать время $T = [0, 1]$.

Пусть $X^{(n)} = (X_t^{(n)}, t \in [0, 1])$, $n \in \mathbb{N}$ - последовательность действительных случайных процессов на $[0, 1]$. Как можно было определить сходимость по распределению такой последовательности?

Идея 1: как слабую сходимость распределений!

Мы поняли, что каждый случайный процесс измерим относительно цилиндрической σ -алгебры на пространстве траекторий. Пусть $F[0, 1]$ - множество вещественных функций на $[0, 1]$, а $\mathcal{C}(F[0, 1])$ - цилиндрическая σ -алгебра на нем. Тогда для процесса $X = (X_t, t \in [0, 1])$ распределение P_X есть вероятностная мера на $(F[0, 1], \mathcal{C}(F[0, 1]))$, определенная по правилу:

$$P_X(C) = P(X \in C) \text{ для любого } C \in \mathcal{C}(F[0, 1]).$$

Замечание

Однако пространство $F[0, 1]$ не является метрическим, в $\mathcal{C}(F[0, 1])$ - не является борелевской σ -алгеброй.

Идея 2: а пусть процесс $X = (X_t, t \in [0, 1])$ имеет непрерывные траектории, тогда его пространством траекторий будет $C[0, 1]$, которое отлично метризуется нормой

$$\|x\| = \max_{t \in [0, 1]} |x(t)|.$$

Замечание

- 1) $C[0, 1]$ не является элементом цилиндрической σ -алгебры $\mathcal{C}(F[0, 1])$;
- 2) X , вообще говоря, не измерим относительно $\mathcal{B}(C[0, 1])$.

Решением этих проблем является рассмотрение цилиндрической σ -алгебры, суженной на $C[0, 1]$:

$$\mathcal{C}(C[0, 1]) = \mathcal{C}(F[0, 1]) \cap C[0, 1] = \{B \cap C[0, 1] : B \in \mathcal{C}(F[0, 1])\}.$$

Лемма

$$\mathcal{B}(C[0, 1]) = \mathcal{C}(C[0, 1]).$$

Доказательство.

Сначала докажем, что $\mathcal{B}(C[0, 1]) \subseteq \mathcal{C}(C[0, 1])$. Для этого возьмем замкнутый шар

$B_r = \{u \in C[0, 1] : \|y - x\| \leq r\}$. Тогда в силу непрерывности всех функций получаем, что

$$B_r(x) = \bigcap_{t \in \mathbb{Q} \cap [0, 1]} \{y \in C[0, 1] : |y(t) - x(t)| \leq r\}.$$

Но это есть ни что иное, как счетное пересечение элементарных цилиндров. Следовательно, $B_r(x) \in \mathcal{C}(C[0, 1])$. Шары порождают всю борелевскую σ -алгебру, поэтому $\mathcal{B}(C[0, 1]) \subseteq \mathcal{C}(C[0, 1])$.

Теперь докажем обратное вложение. Для этого рассмотрим следующий цилиндр: $t \in [0, 1], a < b$,

$$Cyl(t, [a, b]) = \{y \in C[0, 1] : y(t) \in [a, b]\}.$$

Введем функцию $h_t : C[0, 1] \mapsto \mathbb{R}$, действующую по правилу $h_t(x) = x(t)$. Это непрерывная функция, и $Cyl(t, [a, b]) = h_t^{-1}([a, b])$. Отрезок $[a, b]$ замкнут, поэтому $Cyl(t, [a, b])$ является замкнутым множеством в $C[0, 1]$.

Следовательно, $Cyl(t, [a, b]) \in \mathcal{B}(C[0, 1])$. Подобные цилиндры порождают цилиндрическую σ -алгебру, поэтому $\mathcal{C}(C[0, 1]) \subseteq \mathcal{B}(C[0, 1])$. ■

Определение 18.10. Пусть $X^{(n)} = (X_t^{(n)}, t \in [0, 1]), n \in \mathbb{N}$, и $X = (X_t, t \in [0, 1])$ - случайные процессы с непрерывными траекториями. Тогда последовательность $X^{(n)}$ сходится по распределению к X , если их распределения сходятся слабо как вероятностные меры на $C[0, 1]$:

$$P_{X^{(n)}} \xrightarrow{\omega} P_X \text{ при } n \rightarrow \infty.$$

Обозначение: $X^{(n)} \xrightarrow{\mathcal{D}} X$.

Теорема 18.2. о наследовании сходимости

Пусть $h : C[0, 1] \mapsto \mathbb{R}(\mathbb{R}^k)$ - непрерывная функция и $X^{(n)} \xrightarrow{\mathcal{D}} X$ при $n \rightarrow \infty$. Тогда $h(X^{(n)}) \xrightarrow{\mathcal{D}} h(X)$ (как случайные величины или векторы).

Доказательство.

Пусть $f : \mathbb{R}^k \mapsto \mathbb{R}$ - ограниченная непрерывная функция. Тогда композиция $g = f \circ h : C[0, 1] \mapsto \mathbb{R}$ также будет ограниченной непрерывной функцией. Рассмотрим $E f(h(X^{(n)}))$. Оно равно

$$E f(h(X^{(n)})) = E g(X^{(n)}) = \int_{C[0, 1]} g(x) P_{X^{(n)}}(dx).$$

Воспользуемся слабой сходимостью распределений:

$$\lim_{n \rightarrow \infty} \int_{C[0, 1]} g(x) P_{X^{(n)}}(dx) = \int_{C[0, 1]} g(x) P_X(dx) = E g(X).$$

Следовательно, для любой ограниченной непрерывной функции $f : \mathbb{R}^k \mapsto \mathbb{R}$ выполнена следующая сходимость: $E f(h(X^{(n)})) \rightarrow E f(h(X))$ при $n \rightarrow \infty$. Но это означает, что $h(X^{(n)}) \xrightarrow{\mathcal{D}} h(X)$. ■

Принцип инвариантности

Винеровский процесс

Центральным случайным процессом в теории вероятностей является винеровский.

Определение 18.11. Случайный процесс $(W_t, t \geq 0)$ называется винеровским (процессом броуновского движения), если

- 1) $W_0 = 0$ почти наверное;
- 2) W_t имеет независимые приращения, то есть для любого $n \in \mathbb{N}$ и любых $0 \leq t_1 < \dots < t_n$ случайные величины $W_{t_1}, W_{t_2} - W_{t_1}, \dots, W_{t_n} - W_{t_{n-1}}$ независимы в совокупности;
- 3) $W_t - W_s \sim \mathcal{N}(0, t - s), \forall t \geq s \geq 0$.
- 4) Траектории W_t непрерывны почти наверное.

Из свойств винеровского процесса отметим, что W_t является гауссовским, то есть для любого $n \in \mathbb{N}$ и любых $t_1, \dots, t_n \in \mathbb{R}_+$ вектор $(W_{t_1}, W_{t_2}, \dots, W_{t_n})$ является гауссовским.

Ответ на вопрос о том, где возникает сходимость по распределению случайных процессов дает следующая теорема.

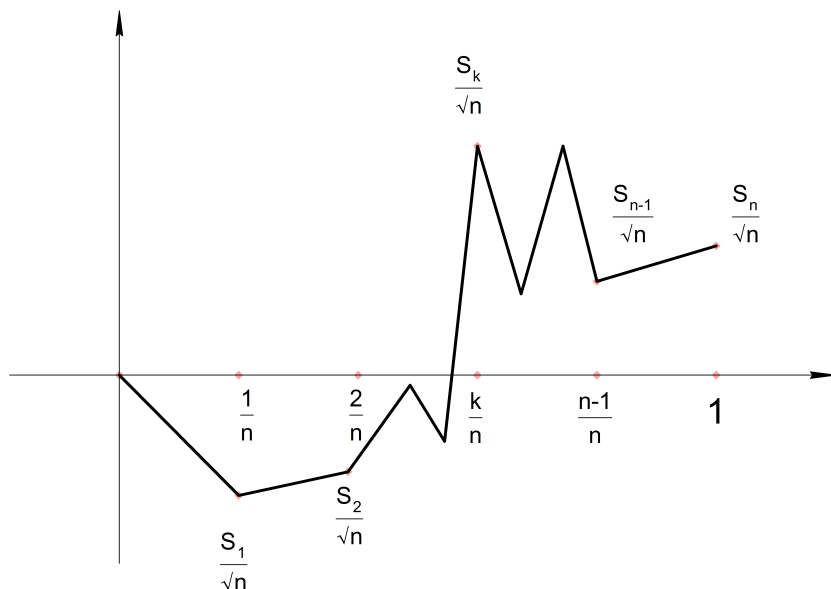


Рис. 7

Теорема 18.3. (принцип инвариантности Донскера-Прохорова)

Пусть $\{\xi_n, n \in \mathbb{N}\}$ - последовательность независимых и одинаково распределенных случайных величин с условием $E\xi_1 = 0$ и $D\xi_1 = 1$. Введем последовательность случайных величин $\{S_n, n \in \mathbb{N}\}$ по следующему правилу:

$X^{(n)} = (X_t^{(n)}, t \in [0, 1])$, как линейную интерполяцию S_0, \dots, S_n :

$$X^{(n)} = \frac{S_k}{\sqrt{n}}(k+1-nt) + \frac{S_{k+1}}{\sqrt{n}}(nt-k) \text{ при } t \in \left[\frac{k}{n}, \frac{k+1}{n}\right], k = 0, \dots, n-1.$$

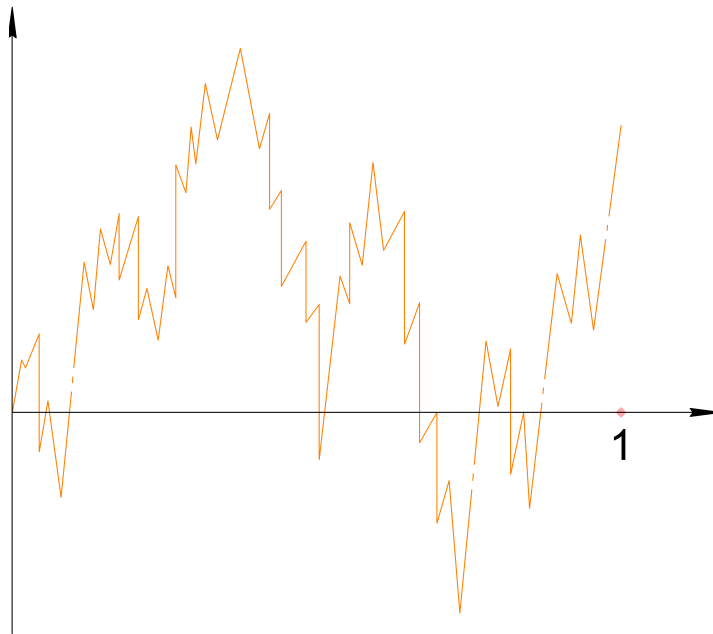


Рис. 8: Винеровский процесс

Тогда $X^{(n)} \xrightarrow{D} W$, где $W = (W_t, t \in [0, 1])$ - винеровский процесс.

Замечание

Центральная предельная теорема является простым следствием принципа инвариантности.

19.

Лекция 19

Теорема Колмогорова

Пусть $(X_n, n \in \mathbb{N})$ - выборка неограниченного размера из распределения с функцией распределения $F(x)$. Построим по ней эмпирическую функцию распределения:

$$\hat{F}_n(x) = \frac{1}{n} \sum_{j=1}^n \mathbf{I}\{X_j \leq x\}.$$

Следующая теорема играет важнейшую роль в статистике для обоснования применения критерия Колмогорова.

Теорема 19.1. (Колмогоров)

Если $F(x)$ непрерывна, то распределение случайной величины

$$D_n = \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)|$$

не зависит от $F(x)$. Более того,

$$\sqrt{n} \cdot D_n \xrightarrow{d} K,$$

где K - распределение Колмогорова с ф.р. $K(y) = \sum_{k \in \mathbb{Z}} (-1)^k e^{-2k^2 y^2}$, $y > 0$.

Комментарий

Если сравнивать теорему Колмогорова с теоремой Гливенко-Кантелли, то здесь аналогия примерно та же, что и сравнение УЗБЧ с ЦПТ. То есть, согласно теореме Гливенко-Кантелли мы знаем, что D_n п.н. сходится к нулю. Теперь мы показываем, с какой скоростью статистика D_n стремится к нулю.

Доказательство теоремы Колмогорова

Разобьем доказательство на несколько шагов.

Лемма (1)

Распределение случайной величины D_n не зависит от $F(x)$.

Доказательство.

Рассмотрим другую выборку: обозначим $Y_j = F(X_j)$. Тогда случайные величины Y_1, \dots, Y_n будут независимыми с распределением $U(0, 1)$: для любого $y \in (0, 1)$

$$P(Y_j \leq y) = P(F(X_j) \leq y) = P(X_j \leq F^{-1}(y)) = F(F^{-1}(y)) = y,$$

где $F^{-1}(y) = \min\{z : F(z) = y\}$. Далее, заметим, что с вероятностью 1

$$I\{X_j \leq x\} = I\{F(X_j) \leq F(x)\} \text{ для всех } x \in \mathbb{R}.$$

Отсюда с вероятностью 1 получаем равенства

$$\begin{aligned} D_n &= \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| = \sup_{x \in \mathbb{R}} \left| \frac{1}{n} \sum_{j=1}^n I\{X_j \leq x\} - F(x) \right| = \\ &= \sup_{x \in \mathbb{R}} \left| \frac{1}{n} \sum_{j=1}^n I\{Y_j \leq F(x)\} - F(x) \right| = \sup_{y \in (0,1)} \left| \frac{1}{n} \sum_{j=1}^n I\{Y_j \leq y\} - y \right|. \end{aligned}$$

Распределение последней случайной величины не зависит от $F(x)$. ■

Лемма (2)

Пусть $(W_t, t \in [0, 1])$ - винеровский процесс. Тогда

$$\sqrt{n}D_n \xrightarrow{d} \max_{t \in [0,1]} |W_t - tW_1|.$$

Доказательство.

В силу леммы (1) можно считать, что $(X_n, n \in \mathbb{N})$ - это выборка из $U(0, 1)$. Поймем, где может достигаться супремум в определении D_n . Раз функция $\hat{F}_n(x)$ является кусочно-постоянной, то супремум разности $|\hat{F}_n(x) - x|$ достигается в точках разрыва. Следовательно,

$$D_n = \max \left(\max_{k=1, \dots, n} \left| X_{(k)} - \frac{k}{n} \right|, \max_{k=1, \dots, n} \left| X_{(k)} - \frac{k-1}{n} \right| \right).$$

Попробуем изобразить то, как получается это равенство. Нарисуем графики эмпирической и настоящей функций распределения. Выборка $(X_1, \dots, X_n) \sim U(0, 1)$, поэтому $F(x) = x$. Эмпирическая функция - кусочно-постоянная функция, которая меняет свое значение в точках порядковых статистик нашей выборки. Значения будут равны: $0, \frac{1}{n}, \frac{2}{n}$, и так далее (см. Рис. 9).

Теперь мы хотим найти значение модуля разности. В точке порядковой статистики значение эмпирической функции распределения равно в точности $\frac{k}{n}$. Соответственно, мы получаем разность $|x_{(k)} - \frac{k}{n}|$, где $x_{(k)}$ - значение настоящей функции распределения. С другой стороны, можно рассмотреть предел слева, тогда разность будет иметь вид $|\frac{k}{n} - x_{(k+1)}|$. Видим, что при некотором k одно из этих двух значений должно давать максимум, так как эмпирическая функция распределения постоянна в окрестности, а настоящая функция $F(x) = x$ растет, поэтому, уменьшая или увеличивая x , мы можем добиваться больших значений по модулю. Отсюда получается такое представление для статистики D_n .

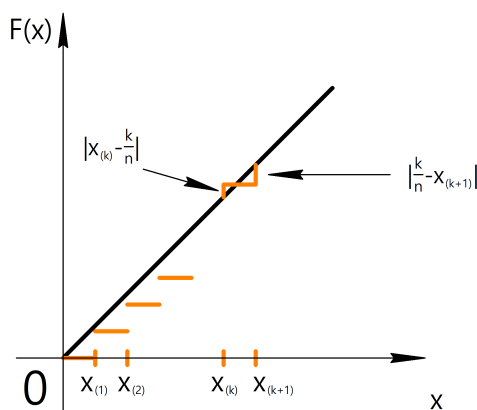


Рис. 9

Обозначим $T_n = \max_{k=1, \dots, n} |X_{(k)} - \frac{k}{n}|$. Тогда $|T_n - D_n| \leq \frac{1}{n}$ и нам достаточно показать, что $\sqrt{n} \cdot T_n \xrightarrow{d} \max_{t \in [0,1]} |W_t - tW_1|$.

Воспользуемся далее следующим упражнением.

Упражнение

Пусть ξ_1, \dots, ξ_{n+1} - это независимые и одинаково распределенные случайные величины с экспоненциальным распределением $\text{Exp}(1)$, а $S_k = \xi_1 + \dots + \xi_k$. Тогда

$$(X_{(1)}, \dots, X_{(n)}) \stackrel{d}{=} \left(\frac{S_1}{S_{n+1}}, \dots, \frac{S_n}{S_{n+1}} \right).$$

Значит, распределение $\sqrt{n} \cdot T_n$ равно распределению

$$\sqrt{n} \cdot \max_{k=1, \dots, n} \left| \frac{S_k}{S_{n+1}} - \frac{k}{n} \right|.$$

Теперь заметим, что в пределе распределение выражения выше будет равно предельному распределению величины

$$\Delta_n = \sqrt{n+1} \cdot \max_{k=1, \dots, n} \left| \frac{S_k}{S_{n+1}} - \frac{k}{n+1} \right|.$$

Преобразуем выражение:

$$\begin{aligned} \Delta_n &= \frac{n+1}{S_{n+1}} \cdot \max_{k=1, \dots, n} \left| \frac{S_k}{\sqrt{n+1}} - \frac{kS_{n+1}}{(n+1)\sqrt{n+1}} \right| = \\ &= \frac{n+1}{S_{n+1}} \cdot \max_{k=1, \dots, n} \left| \frac{S_k - k}{\sqrt{n+1}} - \frac{k}{n+1} \frac{S_{n+1} - (n+1)}{\sqrt{n+1}} \right|. \end{aligned}$$

Теперь введем случайный процесс $Y^{(n)} = (Y_t^{(n)}, t \in [0, 1])$ по следующему правилу: если $t \in [\frac{k}{n+1}, \frac{k+1}{n+1}]$, $k = 0, \dots, n$, то

$$Y_t^{(n)} = \frac{S_k - k}{\sqrt{n+1}}(k+1 - (n+1)t) + \frac{S_{k+1} - (k+1)}{\sqrt{n+1}}((n+1)t - k).$$

Это в точности кусочно-линейный процесс из принципа инвариантности, так как $E\xi_i = 1$ и $D\xi_i = 1$.

Принцип инвариантности говорит нам, что данный процесс будет сходиться по распределению к винеровскому процессу, но нам нужно посмотреть на процесс, который участвует в величине Δ_n : $\frac{S_k - k}{\sqrt{n+1}}$ - это значение нашего процесса в точке $\frac{k}{n+1}$, а $\frac{S_{n+1} - (n+1)}{\sqrt{n+1}}$ - значение в 1. То есть мы рассматриваем процесс $Y_t^{(n)} - tY_1^{(n)}$ и берем его максимальные по модулю значения в точках вида $\frac{k}{n+1}$. Поскольку сам $Y_t^{(n)}$ - это кусочно-линейный процесс, где вершины ломанной - точки вида $\frac{k}{n+1}$, то $Y_t^{(n)} - tY_1^{(n)}$ - это тоже процесс, у которого траектория - ломанная линия, причем вершины этой ломанной будут теми же самыми точками $\frac{k}{n+1}$. Максимальное значение по модулю у траектории, которая является ломанной, достигается в одной из вершин, поэтому

$$\Delta_n = \frac{n+1}{S_{n+1}} \max_{t \in [0,1]} |Y_t^{(n)} - tY_1^{(n)}|.$$

Согласно УЗБЧ выполнено

$$\frac{n+1}{S_{n+1}} \xrightarrow{d} 1.$$

Согласно принципу инвариантности выполнено $Y^{(n)} \xrightarrow{\mathcal{D}} W$, а по теореме о наследовании сходимости получаем, что

$$\max_{t \in [0,1]} |Y_t^{(n)} - tY_1^{(n)}| \xrightarrow{d} \max_{t \in [0,1]} |W_t - tW_1|.$$

Наконец, из леммы Слущкого выводим искомое утверждение:

$$\Delta_n \xrightarrow{d} \max_{t \in [0,1]} |W_t - tW_1|.$$

■

Отметим, что тот процесс, который мы рассматривали в рамках доказательства леммы (2) имеет специальное название.

Определение 19.1. Случайный процесс $W^0 = (W_t^0, t \in [0, 1])$, где $W_t^0 = W_t - tW_1$, называется броуновским мостом.

Мы ответили на основной вопрос: к чему сходится по распределению статистика, которая участвует в формулировке теоремы Колмогорова. Остается вычислить распределение модуля максимума броуновского моста.

Вопрос: как можно вычислить распределение $\max_{t \in [0,1]} |W_t^0|$?

Попробуем понять, как можно по-другому получить распределение броуновского моста, а именно, как условное распределение винеровского процесса при условии, что $W_1 = 0$. Для любого $\varepsilon > 0$ и $B \in \mathcal{B}([0, 1])$ положим

$$Q_\varepsilon(B) = P(W \in B \mid W_1 \in [-\varepsilon, \varepsilon]).$$

Лемма (3)

Пусть P_{W^0} - это распределение броуновского моста. Тогда

$$Q_\varepsilon \xrightarrow{w} P_{W^0} \text{ при } \varepsilon \rightarrow 0+.$$

Доказательство.

По теореме Александрова достаточно проверить, что для любого замкнутого $F \in \mathcal{B}(C[0, 1])$ выполняется неравенство:

$$\overline{\lim}_{\varepsilon \rightarrow 0+} Q_\varepsilon(F) \leq P_{W^0}(F) = P(W^0 \in F).$$

Сначала проверим, что процесс W^0 независим с W_1 . Действительно, для любого $t \in [0, 1]$:

$$\text{cov}(W_t^0, W_1) = \text{cov}(W_t - tW_1, W_1) = 0.$$

Таким образом, для любого $m \in \mathbb{N}$ и для любых $t_1, \dots, t_m \in [0, 1]$ вектор $(W_{t_1}^0, \dots, W_{t_m}^0)$ независим с W_1 . Следовательно, процесс W^0 независим с W_1 .

Далее, рассмотрим расстояние от броуновского моста до винеровского процесса:

$$\|W^0 - W\|_{C[0,1]} = \max_{t \in [0,1]} |W_t - (W_t - tW_1)| = \max_{t \in [0,1]} |tW_1| = |W_1|.$$

Следовательно, если $|W_1| < \delta$ и $W \in F$, то W^0 обязательно принадлежит множеству $F^\delta = \{y \in C[0, 1] : \rho(y, F) \leq \delta\}$. Тогда для любого $\varepsilon < \delta$ выполнено

$$Q_\varepsilon(F) = P(W \in F \mid |W_1| \leq \varepsilon) \leq P(W^0 \in F^\delta \mid |W_1| \leq \varepsilon) = P(W^0 \in F^\delta),$$

где мы в последнем равенстве воспользовались независимостью W^0 и W_1 . Стало быть,

$$\overline{\lim}_{\varepsilon \rightarrow 0+} Q_\varepsilon(F) \leq P(W^0 \in F^\delta) = P_{W^0}(F^\delta).$$

В силу замкнутости F выполнено $F^\delta \downarrow F$ при $\delta \downarrow 0$. По теореме о непрерывности вероятностной меры $P_{W^0}(F^\delta) \rightarrow P_{W^0}(F)$ при $\delta \rightarrow 0$, что и дает искомое неравенство.

■

Следствие (1)

Для любого $z > 0$

$$\begin{aligned} P\left(\max_{t \in [0,1]} |W_t^0| \leq z\right) &= \lim_{\varepsilon \rightarrow 0+} P\left(\max_{t \in [0,1]} |W_t| \leq z \mid W_1 \in [-\varepsilon, \varepsilon]\right) = \\ &= \lim_{\varepsilon \rightarrow 0+} \frac{P\left(\max_{t \in [0,1]} W_t \leq z, \min_{t \in [0,1]} W_t \geq z, W_1 \in [-\varepsilon, \varepsilon]\right)}{P(W_1 \in [-\varepsilon, \varepsilon])}. \end{aligned}$$

Остается найти совместное распределение (m, M, W_1) , где

$$m = \min_{t \in [0,1]} W_t, \quad M = \max_{t \in [0,1]} W_t.$$

Для этого мы решим аналогичную задачу для простейшего симметричного случайного блуждания и применим принцип инвариантности. Пусть $(S_n, n \in \mathbb{N})$ - простейшее симметричное случайное блуждание (S_n представим в виде суммы $\xi_1 + \dots + \xi_n$, где ξ_i - независимые одинаково распределенные случайные величины, которые принимают значение ± 1 с вероятностью $\frac{1}{2}$). Положим

$$m_n = \min_{0 \leq k \leq n} S_k, \quad M_n = \max_{0 \leq k \leq n} S_k.$$

Лемма (4)

Для любых $a \leq 0 \leq b$, $a < b$, $a < v < b$ - целых чисел, выполнено

$$\begin{aligned} P(a < m_n \leq M_n < b, S_n = v) &= \\ &= \sum_{k \in \mathbb{Z}} P(S_n = v + 2k(b-a)) - \sum_{k \in \mathbb{Z}} P(S_n = 2b - v + 2k(b-a)), \end{aligned} \quad (1)$$

где мы считаем, что $S_0 = 0$.

Доказательство.

Будем доказывать индукцией по n . Пусть $n = 0$. Тогда

- 1) Если $v = 0$, $a < 0 < b$, то левая часть равна 1, так как выполняются условия $S_0 = 0$ и $a < 0 \leq 0 < b$. В правой части:
 - Первая сумма равна 1, так как $v = 0$, а $2k(b-a)$ равно нулю только в одном случае: $k = 0$.
 - Вторая сумма равна 0, так как $2b > 0$, $v = 0$, и $2k(b-a) \geq 0$, поэтому S_n не может быть равно 0.
- 2) Если $v \neq 0$, $a < 0 < b$, то левая часть равна 0, так как $S_0 \neq 0$. В правой части:
 - Первая сумма равна 0, так как если $k = 0$, то $S_0 = v \neq 0$. Если же $k > 0$, то $S_0 = v + 2k(b-a) > 0$.

- Вторая сумма равна 0, так как $2b - v > 0$, следовательно $S_0 > 0$.

3) Если a или b равны нулю, то левая часть равна нулю, так как в $v \neq 0$. В правой части обе суммы будут одинаковыми, так как, например при $a = 0$, вероятность во второй сумме - это по сути то же самое, что вероятность того, что $S_n = -v + 2kb$, так как добавление $2b$ - сдвиг индекса k на 1, а S_n - симметричное случайное блуждание, поэтому знак при v не имеет значения в этом случае.

Пусть формула (1) выполнена для $n - 1$, докажем ее для n . Если $a = 0$ или $b = 0$, то снова левая часть равна нулю, а в правой части обе суммы одинаковы.

Так что пусть $a < 0 < b$. Обозначим

$$p_n(a, b, v) = P(a < m_n \leq M_n < b, S_n = v), \quad p_n(j) = P(S_n = j).$$

Тогда заметим, что

$$\begin{aligned} p_n(a, b, v) &= \frac{1}{2}p_{n-1}(a-1, b-1, v-1) + \frac{1}{2}p_{n-1}(a+1, b+1, v+1) \quad \square \\ &\quad (\text{предположение индукции}) \\ &\quad \square \frac{1}{2} \left[\sum_{k \in \mathbb{Z}} p_{n-1}(v-1+2k(b-a)) - \sum_{k \in \mathbb{Z}} p_{n-1}(2b-v+1+2k(b-a)) \right] + \\ &\quad + \frac{1}{2} \left[\sum_{k \in \mathbb{Z}} p_{n-1}(v+1+2k(b-a)) - \sum_{k \in \mathbb{Z}} p_{n-1}(2b-v-1+2k(b-a)) \right] \quad \square \\ &\quad (\text{так как } p_n(j) = \frac{1}{2}p_{n-1}(j-1) + \frac{1}{2}p_{n-1}(j+1)) \\ &\quad \square \sum_{k \in \mathbb{Z}} p_n(v+2k(b-a)) - \sum_{k \in \mathbb{Z}} p_n(2b-v+2k(b-a)). \end{aligned}$$

■

Выведем следствие из леммы (4). Мы хотим расширить диапазон значений S_n , чтобы нам было проще переходить к нормальному закону.

Следствие (2)

Для любых целых чисел $a \leq 0 \leq b$, $a < b$, $a \leq u_1 < u_2 \leq b$ выполнено

$$\begin{aligned} &P(a < m_n \leq M_n < b, u_1 < S_n < u_2) = \\ &= \sum_{k \in \mathbb{Z}} P(u_1 + 2k(b-a) < S_n < u_2 + 2k(b-a)) - \end{aligned}$$

$$- \sum_{k \in \mathbb{Z}} P(2b - u_2 + 2k(b - a) < S_n < 2b - u_1 + 2k(b - a)).$$

Получается суммированием по значениям $v \in (u_1, u_2)$ из леммы (4).

Теперь мы готовы сформулировать лемму о совместном распределении (m, M, W_1) .

Лемма (5)

Для любых вещественных чисел $a < 0 < b$, $a < u < v < b$ выполнено

$$\begin{aligned} & P(a < m \leq M < b, u < W_1 < v) = \\ & = \sum_{k \in \mathbb{Z}} P(u + 2k(b - a) < \xi < v + 2k(b - a)) - \\ & - \sum_{k \in \mathbb{Z}} P(2b - v + 2k(b - a) < \xi < 2b - u + 2k(b - a)), \end{aligned}$$

где $\xi \sim \mathcal{N}(0, 1)$.

Доказательство.

В силу принципа инвариантности и наследования сходимости получаем, что

$$\left(\frac{m_n}{\sqrt{n}}, \frac{M_n}{\sqrt{n}}, \frac{S_n}{\sqrt{n}} \right) \xrightarrow{d} (m, M, W_1).$$

Тогда

$$\begin{aligned} & P(a < m \leq M < b, u < W_1 < v) = \\ & = \lim_{n \rightarrow +\infty} P\left(a < \frac{m_n}{\sqrt{n}} \leq \frac{M_n}{\sqrt{n}} < b, u < \frac{S_n}{\sqrt{n}} < v\right). \end{aligned}$$

Сделаем границы целыми: обозначим $a' = [a\sqrt{n}]$, $b' = [b\sqrt{n}]$, $u' = [u\sqrt{n}]$, $v' = [v\sqrt{n}]$.

Тогда, согласно следствию, получаем:

$$\begin{aligned} & P(a < m \leq M < b, u < W_1 < v) = \\ & = \lim_{n \rightarrow +\infty} \left(\sum_{k \in \mathbb{Z}} P(u' + 2k(b' - a') < S_n < v' + 2k(b' - a')) - \right. \\ & \left. - \sum_{k \in \mathbb{Z}} P(2b' - u' + 2k(b' - a') < S_n < 2b' - v' + 2k(b' - a')) \right). \end{aligned}$$

Заметим, что при фиксированном k в силу центральной предельной теоремы выполнено:

$$\begin{aligned} & \lim_{n \rightarrow +\infty} P(u' + 2k(b' - a') < S_n < v' + 2k(b' - a')) = \\ & = \lim_{n \rightarrow +\infty} P\left(\frac{u' + 2k(b' - a')}{\sqrt{n}} < \frac{S_n}{\sqrt{n}} < \frac{v' + 2k(b' - a')}{\sqrt{n}}\right) = \end{aligned}$$

$$= P(u + 2k(b - a) < \xi < v + 2k(b - a)).$$

Мы хотим переставить предел и суммы местами. Необходимо проверить, что ряды сходятся равномерно по n . Мы знаем, что $\frac{S_n}{\sqrt{n}} \xrightarrow{d} \xi$. Тогда для любого $\varepsilon > 0$ найдется такое $k_0 = k_0(\varepsilon)$, что для любого n выполнено:

$$\begin{aligned} \sum_{k: |k| < k_0} P\left(\frac{u' + 2k(b' - a')}{\sqrt{n}} < \frac{S_n}{\sqrt{n}} < \frac{v' + 2k(b' - a')}{\sqrt{n}}\right) &\leq \\ &\leq P\left(\left|\frac{S_n}{\sqrt{n}}\right| \geq 2k_0(b - a) - |a| - 1\right) \leq \varepsilon. \end{aligned}$$

Переставляя сумму по k и предел по n местами, получаем искомое равенство:

$$\begin{aligned} P(a < m \leq M < b, u < W_1 < v) &= \\ &= \sum_{k \in \mathbb{Z}} P(u + 2k(b - a) < \xi < v + 2k(b - a)) - \\ &- \sum_{k \in \mathbb{Z}} P(2b - v + 2k(b - a) < \xi < 2b - u + 2k(b - a)). \end{aligned}$$

■

Остается подставить правильные параметры в выражение из леммы (5) и посмотреть, куда вероятность при делении на вероятность попадания W_1 в отрезок $[-\varepsilon, \varepsilon]$ (хотим получить функцию распределения Колмогорова).

Найдем распределение случайной величины $\max_{t \in [0,1]} |W_t^0|$. Согласно следствию (1) и лемме (5) для любого $z > 0$

$$\begin{aligned} P\left(\max_{t \in [0,1]} |W_t^0| \leq z\right) &= \lim_{\varepsilon \rightarrow 0+} P\left(\max_{t \in [0,1]} |W_t| \leq z \mid W_1 \in [-\varepsilon, \varepsilon]\right) = \\ &= \lim_{\varepsilon \rightarrow 0+} \left[\frac{1}{P(|W_1| \leq \varepsilon)} \sum_{k \in \mathbb{Z}} P(4kz - \varepsilon < \xi < 4kz + \varepsilon) - \right. \\ &\left. - \frac{1}{P(|W_1| \leq \varepsilon)} \sum_{k \in \mathbb{Z}} P((4k - 2)z - \varepsilon < \xi < (4k - 2)z + \varepsilon) \right]. \end{aligned}$$

Заметим, что

$$\begin{aligned} P(4kz - \varepsilon < \xi < 4kz + \varepsilon) &= \int_{-\varepsilon}^{\varepsilon} \frac{1}{\sqrt{2\pi}} e^{-\frac{(y+4kz)^2}{2}} dy = \\ &= \int_{-\varepsilon}^{\varepsilon} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \cdot e^{-\frac{(4kz)^2}{2} - (4kz)y} dy \leq e^{-\frac{(4kz)^2}{2} + |(4kz)|\varepsilon} \cdot P(|W_1| \leq \varepsilon). \end{aligned}$$

Ряд $\sum_{k \in \mathbb{Z}} e^{-\frac{(4kz)^2}{2} + |(4kz)|\varepsilon}$ сходится равномерно по $\varepsilon \leq 1$. Значит можно переставить местами сумму по k и предел по ε .

Для фиксированного $k \in \mathbb{Z}$ и $\varepsilon \rightarrow 0$ имеем:

$$P(4kz - \varepsilon < \xi < 4kz + \varepsilon) \sim \frac{1}{\sqrt{2\pi}} \cdot 2\varepsilon \cdot e^{-\frac{(4kz)^2}{2}}, \text{ а } P(|W_1| \leq \varepsilon) \sim \frac{1}{\sqrt{2\pi}} \cdot 2\varepsilon.$$

В итоге,

$$\begin{aligned} P\left(\max_{t \in [0,1]} |W_t^0| \leq z\right) &= \sum_{k \in \mathbb{Z}} e^{-\frac{(4kz)^2}{2}} - \sum_{k \in \mathbb{Z}} e^{-\frac{(4k-2)z)^2}{2}} = \\ &= \sum_{k \in \mathbb{Z}} (-1)^k e^{-2k^2 z^2} = K(z). \end{aligned}$$

Доказательство теоремы Колмогорова завершено.

20.

Лекция 20

Это последняя лекция курса Математической статистики, на которой мы познакомимся с элементами последовательного анализа - разделом математической статистики, который тесно связан с теорией случайных процессов.

Мотивировка

Вернемся к задаче, с которой мы начинали рассматривать проверку гипотез. Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из неизвестного распределения P , о котором высказаны две простые гипотезы: $H_0 : P = P_0$ и альтернатива $H_1 : P = P_1$. Относительно распределений P_0 и P_1 известно, что они принадлежат одному доминируемому семейству и имеют плотности $p_0(x)$ и $p_1(x)$.

Казалось бы, все уже изучено. Согласно лемме Неймана-Пирсона р.н.м.к. уровня значимости ε_0 задается видом $S_\lambda = \{\mathbf{X} : p_1(\mathbf{X}) \geq \lambda p_0(\mathbf{X})\}$, а параметра λ находится из уравнения

$$P_0(\mathbf{X} \in S_\lambda) = \varepsilon_0.$$

Но есть одна неприятность, которая состоит в том, что мы мало что знаем о вероятности ошибки второго рода. Вероятность ошибки первого рода равна ε_0 , однако вероятность ошибки второго рода $\varepsilon_1 = P_1(\mathbf{X} \notin S_\lambda)$, хоть и имеет минимально возможное значение, может быть достаточно большой. Но в случае состоятельности критерия вероятность ошибки второго рода будет мала при больших n .

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из схемы Бернулли $\text{Bin}(1, \theta)$. Тогда р.н.м.к. уровня значимости $\varepsilon_0 = 0.05$ для проверки гипотезы $H_0 : \theta = 0.05$ против альтернативы $H_1 : \theta = 0.17$ будет иметь вероятность ошибки второго рода $\varepsilon_1 \leq 0.1$ начиная с $n = 57$.

Но наблюдения могут быть дороги в получении. С практической точки зрения задача минимизации затрат весьма принципиальна. Для ее решения был предложен так называемый метод последовательного анализа, в рамках которого данные поступают по одному, и мы на каждом шаге должны либо принять одну из гипотез, либо потребовать следующее наблюдение.

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка, по которой мы хотим различить две гипотезы

$H_0 : X_i \sim U(1, 3)$ против альтернативы $H_1 : X_i \sim U(0, 2)$. Но если уже X_1 попадает в $[2, 3]$ или $[0, 1]$, то мы можем однозначно сказать, какая из гипотез верна. Остальные данные вообще не потребуются.

Последовательные критерии

Итак, пусть P - неизвестное распределение элементов выборки, о котором высказаны две гипотезы: основная $H_0 : P = P_0$ и альтернатива $H_1 : P = P_1$. Пусть \mathcal{X} - выборочное пространство для одного элемента выборки.

Определение 20.1. Последовательным критерием S для проверки H_0 против H_1 называется такая последовательность пар $S = ((S_0^n, S_1^n), n \in \mathbb{N})$, что для любого $n \in \mathbb{N}$,

$$S_0^n, S_1^n \subset \mathcal{X}^n \text{ и } S_0^n \cap S_1^n = \emptyset.$$

При этом правильно принятия гипотез состоит в следующем. Пусть

$$\tau(S) = \min\{n : \mathbf{X}_n = (X_1, \dots, X_n) \in S_0^n \sqcup S_1^n\},$$

тогда гипотеза H_1 принимается, если $\mathbf{X}_{\tau(S)} = (X_1, \dots, X_{\tau(S)}) \in S_1^{\tau(S)}$, $i = 0, 1$.

Замечание

Обычные критерии являются частным случаем последовательных.

Изобразим разницу между обычными и последовательными критериями.

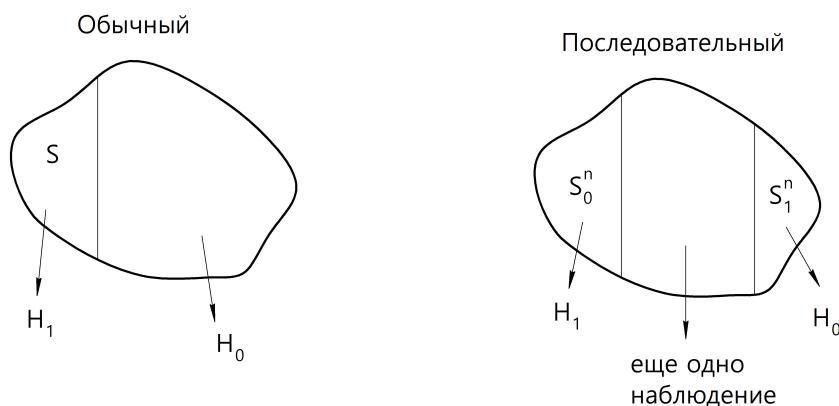


Рис. 10

В обычном критерии мы делили пространство на две части: множество S (критерий), попадая в который, мы принимаем гипотезу H_1 , если мы попадаем не в S , то мы принимаем гипотезу H_0 .

В последовательном критерии выборочное пространство делится на 3 части: S_0^n , S_1^n и дополнение к их объединению. Попадаем в S_0^n - принимаем H_0 , попадаем в S_1^n - принимаем H_1 , в противном случае мы требуем еще одно наблюдение.

Вероятности ошибок первого и второго рода для последовательного критерия равны

$$\alpha_0(S) = P_0(\mathbf{X}_{\tau(S)} \in S_1^{\tau(S)}), \quad \alpha_1(S) = P_1(\mathbf{X}_{\tau(S)} \in S_0^{\tau(S)})$$

Критерий отношения правдоподобия

Задача

Для заданных $\varepsilon_0, \varepsilon_1 \in (0, 1)$ найти последовательный критерий S с условиями

$$\alpha_0(S) = P_0(\mathbf{X}_{\tau(S)} \in S_1^{\tau(S)}) \leq \varepsilon_0, \quad \alpha_1(S) = P_1(\mathbf{X}_{\tau(S)} \in S_0^{\tau(S)}) \leq \varepsilon_1,$$

который одновременно минимизирует величины $E_i \tau(S)$, $i = 0, 1$.

Для решения задачи оптимальным будет критерий, который следует логике леммы Неймана-Пирсона и использует отношение правдоподобия. Пусть, как и ранее, распределение P_0 и P_1 взяты из одного доминируемого семейства и имеют плотности $p_0(x)$ и $p_1(x)$. Введем плотность выборки:

$$p_i(X_1, \dots, X_n) = \prod_{j=1}^n p_i(X_j), \quad i = 0, 1.$$

Определение 20.2. Пусть $0 < \Gamma_0 < 1 < \Gamma_1$ - некоторые числа, $\Gamma = (\Gamma_0, \Gamma_1)$. Последовательным критерием отношения правдоподобия $S(\Gamma)$ называется последовательный критерий, у которого множества S_0^n и S_1^n определены следующим образом:

$$S_0^n = \left\{ \mathbf{X}_n : \frac{p_1(X_1, \dots, X_n)}{p_0(X_1, \dots, X_n)} \leq \Gamma_0 \right\},$$

$$S_1^n = \left\{ \mathbf{X}_n : \frac{p_1(X_1, \dots, X_n)}{p_0(X_1, \dots, X_n)} \geq \Gamma_1 \right\}.$$

Байесовский подход к проверке гипотез

Мы хотим доказать, что для некоторого Γ такой критерий будет оптимальным последовательным критерием. Однако, перед этим нам понадобится установить, что такой же критерий будет оптимальным в байесовском подходе к последовательной проверке гипотез. Пусть у нас есть положительные числа w_0, w_1, a , которые имеют смысл "штрафов" :

- 1) w_i - штраф за совершение ошибки $(i + 1)$ -го рода, $i = 0, 1$;
- 2) a - штраф за взятие дополнительного наблюдения.

Пусть далее, у нас есть априорное распределение вероятностей на множестве гипотез. В нашем случае - это просто число $q \in [0, 1]$, так как гипотез всего две.

Определение 20.3. Байесовским риском критерия S называется величина

$$R(q, S) = q[\alpha_0(S)w_0 + aE_0\tau(S)] + (1 - q)[\alpha_1(S)w_1 + aE_1\tau(S)].$$

Последовательный критерий, который минимизирует (для заданных q, w_0, w_1, a) байесовский риск называется байесовским.

Оптимальность в байесовском смысле

Теорема 20.1. Для любых (q, w_0, w_1, a) существуют такие Γ_1 и Γ_2 , что критерий $S(\Gamma)$ является байесовским.

Доказательство.

Обозначим через δ_i критерий, который принимает гипотезу H_i , $i = 0, 1$, без проведения наблюдений. Тогда $\tau(\delta_i) = 0$, $\alpha_i(\delta_i) = 0$. Поймем, когда критерий S , минимизирующий байесовский риск $R(q, S)$, совпадает с δ_i . Очевидно, что

$$R(q, \delta_0) = (1 - q)w_1, \quad R(q, \delta_1) = qw_0.$$

Пусть \mathcal{K} - класс критериев, которые зависят хотя бы от одного наблюдения. Тогда $R(q, S) \geq a$ для любого $S \in \mathcal{K}$. Положим

$$R(q) = \inf_{S \in \mathcal{K}} R(q, S).$$

Ясно, что $R(q) < +\infty$, так как величина конечна для любого критерия, который требует ровно одно наблюдение.

Проверим, что $R(q)$ - выпуклая кверху функция. В силу линейности (по q) величины $R(q, S)$, получаем следующее: для любого $p \in (0, 1)$

$$\begin{aligned} R(pq_1 + (1 - p)q_2, S) &= \\ &= (pq_1 + (1 - p)q_2)[\alpha_0(S)w_0 + aE_0\tau(S)] + (1 - pq_1 - (1 - p)q_2)[\alpha_1(S)w_1 + aE_1\tau(S)] \\ &= pR(q_1, S) + (1 - p)R(q_2, S). \end{aligned}$$

Следовательно,

$$R(pq_1 + (1 - p)q_2) = \inf_{S \in \mathcal{K}} [pR(q_1, S) + (1 - p)R(q_2, S)] \geq pR(q_1) + (1 - p)R(q_2).$$

Значит, функция $R(q)$ выпукла кверху. И так как $a \leq R(q) < +\infty$ мы получаем, что $R(q)$ непрерывна на $[0, 1]$. Сравним теперь риски критериев δ_0 , δ_1 и $S \in \mathcal{K}$.

Критерию δ_0 соответствует функция байесовского риска $(1 - q)w_1$, критерию δ_1 соответствует qw_0 , и есть $R(q) \geq a$ (значение байесовского риска, которое дает оптимальный критерий для заданного q).

Тогда на отрезке $[0, 1]$ δ_0 - убывающая функция, которая равна нулю в 1. δ_1 - наоборот, возрастающая функция.

Дальше все зависит от того, как пройдет функция $R(q)$. Она ограничена снизу значением a .

- 1) Если точка пересечения прямых δ_0 и δ_1 попала внутрь $R(q)$, то наблюдений брать не нужно, так как точка пересечения - это точка $\frac{w_1}{w_0 + w_1}$, в которой минимален максимальный байесовский риск из пары δ_0, δ_1 . Соответственно, если в этой точке значение функции $R(q)$ больше, чем в точке пересечения, то есть

$$R\left(\frac{w_1}{w_0 + w_1}\right) \geq \frac{w_0 w_1}{w_0 + w_1},$$

то в силу выпуклости кверху функция $R(q)$ будет выше, чем минимум из пары рисков δ_0, δ_1 . Это значит, что мы всегда предпочитаем первую или вторую гипотезу в зависимости от q : если $q \in [0, \frac{w_1}{w_0 + w_1})$, то мы предпочитаем H_1 , если $q \in (\frac{w_1}{w_0 + w_1}, 1]$, то H_0 .

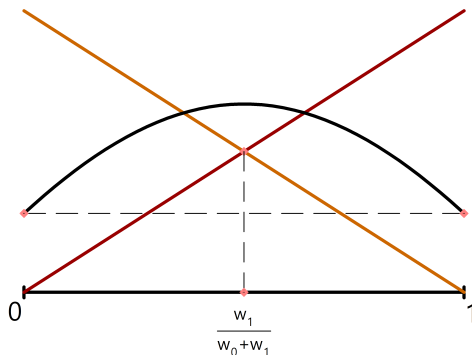


Рис. 11

- 2) Пусть теперь

$$R\left(\frac{w_1}{w_0 + w_1}\right) < \frac{w_0 w_1}{w_0 + w_1}.$$

Мы получаем интервал значений q , в котором нам предпочтительней взять ещё одно наблюдение, чем пытаться выбрать кого-то из пары δ_0, δ_1 . Этот интервал ограничен точками, в которых $R(q)$ пересекает прямые. Тогда, попадая левее или правее от этого интервала, мы предпочитаем, соответственно, δ_1 или

δ_0 . В интервале оптимальным будет критерий, имеющий хотя-бы одно наблюдение.

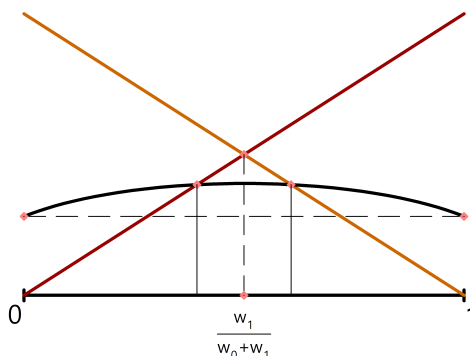


Рис. 12

Более формально эта ситуация выглядит так: существуют решения уравнений $R(q, \delta_i) = R(q)$, $i = 0, 1$. Обозначим их через $1 - \gamma_0 > 1 - \gamma_1$. В первом случае положим

$$1 - \gamma_0 = 1 - \gamma_1 = \frac{w_1}{w_0 + w_1}.$$

Тогда из приведенных рассуждений вытекает следующий оптимальный план действий. По заданным w_0, w_1, a мы вычисляем величины $1 - \gamma_0, 1 - \gamma_1$.

- Если $q \leq 1 - \gamma_1$ (или $\gamma_1 \leq 1 - q$), то наименьший риск имеет критерий δ_1 и надо принять H_1 .
- Если $q \geq 1 - \gamma_0$ (или $\gamma_0 \geq 1 - q$), то наименьший риск имеет критерий δ_0 и надо принять H_0 .
- Если же $q \in (1 - \gamma_1, 1 - \gamma_0)$ (или $1 - q \in (\gamma_0, \gamma_1)$), то оптимален кто-то из \mathcal{K} , то есть надо провести наблюдение.

Проведем теперь индукцию по n . Пусть проведено n наблюдений и имеем выборку (X_1, \dots, X_n) . Перед $(n + 1)$ -м наблюдением мы имеем ту же альтернативу: либо не проводить больше наблюдений и принять одну из гипотез, либо провести еще одно наблюдение. Тот факт, что риск уже увеличен на an , роли не играет, так как эти штрафные потери уже не могут быть устранены. Однако, теперь у нас будут другие априорные вероятности для гипотез.

По факту они превращаются в апостериорные вероятности параметра гипотезы при условии заданных X_1, \dots, X_n :

$$q(0|\mathbf{X}_n) = \frac{qp_0(X_1, \dots, X_n)}{qp_0(X_1, \dots, X_n) + (1 - q)p_1(X_1, \dots, X_n)},$$

$$q(1|\mathbf{X}_n) = \frac{(1-q)p_1(X_1, \dots, X_n)}{qp_0(X_1, \dots, X_n) + (1-q)p_1(X_1, \dots, X_n)}.$$

Далее, мы уже знаем оптимальное решение:

- Если $q(0|\mathbf{X}_n) \leq 1 - \gamma_1$, то надо принять H_1 .
- Если $q(0|\mathbf{X}_n) \geq 1 - \gamma_0$, то надо принять H_0 .
- Если же $q(0|\mathbf{X}_n) \in (1 - \gamma_1, 1 - \gamma_0)$ (или $q(1|\mathbf{X}_n) \in (\gamma_0, \gamma_1)$), то надо провести наблюдение.

Таким образом, мы продолжаем наблюдения, если

$$\gamma_0 < \frac{(1-q)p_1(X_1, \dots, X_n)}{qp_0(X_1, \dots, X_n) + (1-q)p_1(X_1, \dots, X_n)} < \gamma_1$$

Перепишем эти неравенства в терминах ограничений на отношение p_1 к p_0

$$\frac{q\gamma_0}{(1-q)(1-\gamma_0)} < \frac{p_1(X_1, \dots, X_n)}{p_0(X_1, \dots, X_n)} < \frac{q\gamma_1}{(1-q)(1-\gamma_1)}.$$

Но это и есть в точности критерий $S(\Gamma)$ для $\Gamma_i = \frac{q\gamma_i}{(1-q)(1-\gamma_i)}$, $i = 0, 1$. ■

Замечание

Заметим, что числа $\gamma_i = \gamma_i(w_0, w_1, a)$, $i = 0, 1$ остаются теми же, если все параметры будут умножены на одно и то же положительное число.

Тем самым, у нас есть лишь два параметра. Для удобства будем считать, что $w_1 = 1 - w_0$. Важна следующая техническая лемма, которую, впрочем, мы оставим без доказательства.

Лемма (1)

Для любых $0 < \gamma_0 < \gamma_1 < 1$ найдутся такие $w \in (0, 1)$ и $a > 0$, что для любого $q \in (1 - \gamma_1, 1 - \gamma_0)$ оптимальное решение байесовской задачи с параметрами $w_0 = 1 - w$, $w_1 = w$, a и априорной вероятностью q задается последовательным критерием отношением правдоподобия $S(\Gamma)$, где

$$\Gamma_0 = \frac{q\gamma_0}{(1-q)(1-\gamma_0)}, \quad \Gamma_1 = \frac{q\gamma_1}{(1-q)(1-\gamma_1)}.$$

Оптимальность последовательного критерия

Теперь мы готовы сформулировать и доказать теорему об оптимальности последовательного критерия отношения правдоподобия.

Теорема 20.2. (об оптимальности)

Пусть $\Gamma_0 < 1 < \Gamma_1$. Обозначим через ε_i , $i = 0, 1$, вероятности ошибок первого и второго родов последовательного критерия отношения правдоподобия $S(\Gamma)$. Тогда среди всех последовательных критериев S , для которых

$$\alpha_i(S) \leq \varepsilon_i, \quad i = 0, 1,$$

критерий $S(\Gamma)$ имеет наименьшие значения величин $E_i \tau(S)$, $i = 0, 1$.

Доказательство.

Сведем задачу к байесовскому случаю. Возьмем произвольное $q \in (0, 1)$ и найдем величины γ_0, γ_1 из равенств:

$$\Gamma_i = \frac{q\gamma_i}{(1-q)(1-\gamma_i)}, \quad i = 0, 1.$$

Заметим, что в силу условия $\Gamma_0 < 1 < \Gamma_1$ мы автоматически получим, что

$$q \in (1 - \gamma_1, 1 - \gamma_0).$$

Тогда согласно лемме (1) найдутся такие параметры a, w , что критерий $S(\Gamma)$ будет байесовским в задаче с параметрами $(q, 1 - w, w, a)$.

Пусть S - теперь любой другой критерий с условиями

$$\alpha_i(S) \leq \varepsilon_i, \quad i = 0, 1.$$

Сравним теперь S и $S(\Gamma)$ в байесовском подходе с параметрами $(q, 1 - w, w, a)$. Мы знаем, что критерий $S(\Gamma)$ минимизирует байесовский риск, поэтому выполняется неравенство

$$\begin{aligned} & q[\alpha_0(S)w_0 + aE_0\tau(S)] + (1-q)[\alpha_1(S)w_1 + aE_1\tau(S)] \geq \\ & \geq q[\alpha_0(S(\Gamma))w_0 + aE_0\tau(S(\Gamma))] + (1-q)[\alpha_1(S(\Gamma))w_1 + aE_1\tau(S(\Gamma))]. \end{aligned}$$

Из этого неравенства и условий на S следует, что

$$qE_0\tau(S) + (1-q)E_1\tau(S) \geq qE_0\tau(S(\Gamma)) + (1-q)E_1\tau(S(\Gamma)).$$

В силу произвольности $q \in (0, 1)$ мы получаем искомые соотношения:

$$E_0\tau(S) \geq E_0\tau(S(\Gamma)), \quad E_1\tau(S) \geq E_1(\tau(S(\Gamma))).$$

■

Приведем важное замечание.

Замечание

Можно показать, что если распределения P_0 и P_1 непрерывны, то для любых $\varepsilon_0, \varepsilon_1 \in (0, 1)$ найдутся такие $\Gamma_0 < 1 < \Gamma_1$, что последовательный критерий $S(\Gamma)$ будет иметь вероятности ошибок первого и второго рода в точности ε_0 и ε_1 .

К сожалению, явное нахождение величин $\Gamma_0 < 1 < \Gamma_1$ при заданных вероятностях ошибок весьма затруднительно. По сути, оно представляет собой явное нахождение границ (Γ_0, Γ_1) , для которых уже заданы вероятности первого выхода случайного блуждания из полосы. Даже для схемы Бернулли необходима правильная связь между распределениями.

Упражнение

Рассмотрим задачу построения критерия последовательного отношения правдоподобия для выборки из биномиального распределения $\text{Bin}(1, \theta)$, $H_0 : \theta = p_0 < \frac{1}{2}$, $H_1 : \theta = p_1 > \frac{1}{2}$. Положим $q_i = 1 - p_i$. Предположим, что $p_1 = q_0$ и числа Γ_0, Γ_1 подобраны так, что

$$\frac{\ln \Gamma_0}{\ln\left(\frac{q_0}{p_0}\right)} = -a, \quad \frac{\ln \Gamma_1}{\ln\left(\frac{q_0}{p_0}\right)} = b,$$

где a и b - положительные целые числа. Вычислите вероятности ошибок первого и второго рода для критерия $S(\Gamma)$.

Приближенные вычисления

Если мы не можем решить задачу точно, то попробуем найти близкие вероятности ошибок, для которых уже можно вычислить границы Γ . Например, можно пойти по следующему пути.

Введем события:

$$A_n = \left\{ \Gamma_0 < \frac{p_1(X_1, \dots, X_k)}{p_0(X_1, \dots, X_k)} < \Gamma_1, \quad k = 1, \dots, n-1, \quad \frac{p_1(X_1, \dots, X_k)}{p_0(X_1, \dots, X_k)} \leq \Gamma_0 \right\};$$

$$B_n = \left\{ \Gamma_0 < \frac{p_1(X_1, \dots, X_k)}{p_0(X_1, \dots, X_k)} < \Gamma_1, \quad k = 1, \dots, n-1, \quad \frac{p_1(X_1, \dots, X_k)}{p_0(X_1, \dots, X_k)} \geq \Gamma_1 \right\}.$$

Ясно, что эти события не пересекаются. $\cup_n A_n$ есть событие, что мы принимаем H_0 , а $\cup_n B_n$ есть событие, что мы принимаем H_1 . Тогда

$$\begin{aligned} \alpha_0(S(\Gamma)) &= \sum_{n=1}^{+\infty} P_0(B_n) = \sum_{n=1}^{+\infty} \int_{\mathcal{X}^n} I_{B_n}(x_1, \dots, x_n) p_0(x_1, \dots, x_n) \mu(d\mathbf{x}) \leq \\ &= \frac{1}{\Gamma_1} \sum_{n=1}^{+\infty} \int_{\mathcal{X}^n} I_{B_n}(x_1, \dots, x_n) p_1(x_1, \dots, x_n) \mu(d\mathbf{x}) = \\ &= \frac{1}{\Gamma_1} \sum_{n=1}^{+\infty} P_1(B_n) = \frac{1 - \alpha_1(S(\Gamma))}{\Gamma_1}. \end{aligned}$$

Совершенно аналогично получаем, что

$$\alpha_1(S(\Gamma)) \leq \Gamma_0(1 - \alpha_0(S(\Gamma))).$$

Обозначим для краткости $\alpha_i = \alpha_i(S(\Gamma))$, $i = 0, 1$. Тогда

$$\Gamma_0 \geq \frac{\alpha_1}{1 - \alpha_0} = \Gamma'_0, \quad \Gamma_1 \leq \frac{1 - \alpha_1}{\alpha_0} = \Gamma'_1. \quad (1)$$

Мы получили некоторые оценки на границы, которые мы переобозначили штрихами. Рассмотрим теперь критерий $S(\Gamma')$. Его вероятности ошибок α'_i , $i = 0, 1$ в силу (1) будут удовлетворять следующим соотношениям:

$$\frac{\alpha_1}{1 - \alpha_0} \geq \frac{\alpha'_1}{1 - \alpha'_0}, \quad \frac{1 - \alpha_1}{\alpha_0} \leq \frac{1 - \alpha'_1}{\alpha'_0}.$$

Следовательно

$$\begin{aligned} \alpha'_0 &\leq \frac{\alpha_0(1 - \alpha'_1)}{1 - \alpha_1} \leq \frac{\alpha_0}{1 - \alpha_1}, \quad \alpha'_1 \leq \frac{\alpha_1(1 - \alpha'_0)}{1 - \alpha_0} \leq \frac{\alpha_1}{1 - \alpha_0} \\ \alpha'_0 + \alpha'_1 &\leq \alpha_0 + \alpha_1. \end{aligned}$$

Таким образом, для нового критерия $S(\Gamma')$ вероятности ошибок возможно будут лишь немного больше (при малых α_0, α_1), а сумма вероятностей ошибок - даже меньше.

Численный пример

Вернемся к примеру из начала лекции.

Пример

Пусть $\mathbf{X} = (X_1, \dots, X_n)$ - выборка из схемы Бернулли $(1, \theta)$. Тогда р.н.м.к. уровня значимости $\varepsilon_0 = 0.05$ для проверки гипотезы $H_0 : \theta = 0.05$ против альтернативы $H_1 : \theta = 0.17$ будет иметь вероятность ошибки второго рода $\varepsilon_1 \leq 0.1$ начиная с $n = 57$.

С другой стороны, для данного примера при $\alpha_0 = 0.05$, $\alpha_1 = 0.1$ будет выполнено

$$\Gamma'_0 = 0.105, \quad \Gamma'_1 = 18, \quad \alpha'_0 = 0.031, \quad \alpha'_1 = 0.099,$$

но

$$E_0 \tau(S(\Gamma')) = 31.4, \quad E_1 \tau(S(\Gamma')) = 30.$$

Таким образом, последовательный критерий требует в среднем почти в два раза меньшее число наблюдений при почти тех же вероятностях ошибок первого и второго рода.

Тем самым, мы показали, что последовательный критерий действительно имеет преимущества: мы значительно выигрываем в количестве наблюдений, что крайне удобно с практической точки зрения.



МЕХАНИКО-
МАТЕМАТИЧЕСКИЙ
ФАКУЛЬТЕТ
МГУ ИМЕНИ
М.В. ЛОМОНОСОВА

teach-in
Л Е К Ц И И У Ч Е Н Ы Х М Г У