

Лекции по Математической Статистике 5 семестр

Луа Yaroshevskiy

13 мая 2023 г.

Оглавление

Лекция 1	2
1.1 Введение	2
1.2 Выборочные характеристики	2
1.3 Первоначальная обработка	3
Лекция 2	5
2.1 Свойства статистических оценок	5
2.2 Точечные оценки моментов	6
Лекция 3	8
3.1 Метод максимального правдоподобия	8
3.2 Неравенство Раво-Крамера	9
TODO Лекция 4	11
TODO Лекция 5	12
TODO Лекция 6	13
TODO Лекция 7	14
Лекция 9	15
8.1 Математическая регрессия	15
8.1.1 Метод наименьшие квадратов	15
8.2 Линейная парная регрессия	16
8.2.1 Геометрический смысл прямой линейной регрессии	16
8.2.2 Проверка гипотезы по значимости коэффициента линейной корреляции	17
8.2.3 Выборочное корреляционное отношение	17

Лекция 1

1.1 Введение

Определение. Генеральная совокупность — множество всех исходов определенного всех экспериментов

Определение. Выборочная совокупность — множество исходов наблюдаемых экспериментов

Примечание. Выборка репрезентативная если ее распределение совпадает с распределением генеральной совокупности

Определение. 1 Пусть проведено n наблюдаемых независимых экспериментов в которых наблюдаемые величины приняли значения: X_1, X_2, \dots, X_n . Набор этих данных называется **выборкой объема n**

Определение. 2 **Выборкой объема n** называется набор из n независимых одинаково распределенных случайных величин

1.2 Выборочные характеристики

Пусть имеется выборка в смысле 1. Ее можно понимать как следующую дискретную случайную величину:

$$\frac{X_1}{p_i} \quad \frac{X_2}{\frac{1}{n}} \quad \frac{X_3}{\frac{1}{n}} \quad \dots \quad \frac{X_n}{\frac{1}{n}} \quad \frac{\sum}{1}$$

$$X = \sum_{i=1}^n \frac{1}{n} X_i = \frac{1}{n} \sum_{i=1}^n X_i$$

— матожидание?

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

— дисперсия

$$F_n^*(z) = \frac{1}{n} \sum_{i=1}^n I(X_i < z)$$

— функция распределения, где $I(X < z) = \begin{cases} 1 & , X < z \\ 0 & , X \geq z \end{cases}$ — индикатор

Теорема 1.2.1. $\forall z \in \mathbb{R}$

$$F_n^*(z) \xrightarrow[n \rightarrow \infty]{p} F(z)$$

Доказательство.

$$EI(X_1 < z) = 1 \cdot p(X < z) + 0 \cdot p(X_1 \geq z) = p(X_1 < z) = F(z)$$

— функция распределения X_1

По ЗБЧ Хнчина:

$$F_n^*(z) = \frac{\sum I(X_i < z)}{n} \xrightarrow{p} EI(X_1 < z) = F(z)$$

□

Примечание. На самом деле имеется даже равномерная сходимость по вероятности: *теорема Гливленко-Кантеля:*

$$\sup_{z \in \mathbb{R}} |F_n^*(z) - F(z)| \xrightarrow{p} [n \rightarrow \infty] 0$$

1.3 Первоначальная обработка

Определение. Если упорядочить выборку по возрастанию (ранжирование) то получим **вариационный ряд**:

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$

Если учесть повторяющиеся экземпляры, получим **частотный вариационный ряд**

$$\begin{array}{cccccc} X_{(1)} & X_{(2)} & \dots & X_{(k)} & \sum & \\ \hline n_1 & n_2 & \dots & n_k & n & \\ \hline p_1 & \frac{n_1}{n} & \dots & \frac{n_k}{n} & 1 & \end{array}$$

- $h = X_{\max} - X_{\min}$ — **размах выборки**

Примечание. Допустим разбили интервал (X_{\min}, X_{\max}) на k интервалов, чаще всего одинаковой длины. Тогда $h = \frac{h}{k}$. Тогда вариационный ряд можно заменить интервальным вариационным рядом.

Пример.

$$\begin{array}{cccccc} i & l_1 & l_2 & \dots & l_k & \sum \\ \hline m_i & m_2 & m_2 & \dots & m_k & n \\ \hline \frac{m_i}{n} & \frac{m_1}{m} & \frac{m_2}{m} & \dots & \frac{m_k}{m} & 1 \end{array}$$

На координатной плоскости построим прямоугольники:

- l_i — основание прямоугольника соответствующего интервала
- $\frac{m_i}{nl_i}$ — высота прямоугольника

Получаем ступенчатую фигуру площади 1, которая называется **гистограмма**.

Теорема 1.3.1. При

- $n \rightarrow \infty$

- $k(n) \rightarrow \infty$
- $\frac{k(n)}{n} \rightarrow 0$

Гистограмма по вероятности будет стремиться к плотности распределения:

$$\frac{m_i}{n} \xrightarrow{p} p(X_i \in l_i) = \int_{l_i} f(x) dx$$

Примечание. Чаще всего число интервалов берется по формулу Стержеса.

$$k \approx 1 + \log_k n$$

$$k \approx \sqrt[3]{n}$$

Примечание. Иногда выборка изображается в виде полигона

- l_i^m — середина i -го интервала

Лекция 2

$$\mathbb{X} = (\mathbb{X}_1, \dots, \mathbb{X}_n)$$

Определение. Статистикой называется измеримая функция:

$$\Theta^* = \Theta^*(\mathbb{X}_1, \dots, \mathbb{X}_n)$$

Примечание. Пусть требуется найти оценку параметра Θ случайной величины \mathbb{X} по данной выборке. Оценку будет считать с помощью некоторой статистики:

$$\Theta^* = \Theta^*(\mathbb{X}_1, \dots, \mathbb{X}_n)$$

2.1 Свойства статистических оценок

1. Состоятельность

Определение. Статистика Θ^* называется **состоятельной оценкой** параметра Θ , если

$$\Theta^* \xrightarrow{P} \Theta \text{ при } n \rightarrow \infty$$

2. Несмещенность

Определение. Статистика Θ^* называется несмещенной оценкой параметра Θ , если

$$E\Theta^* = \Theta$$

Примечание. С равной вероятностью можем ошибиться как в меньшую так и в большую сторону. Нет систематической ошибки.

Определение. Статистика Θ^* называется **асимптотически несмещенной**, если

$$\Theta^* \xrightarrow[n \rightarrow \infty]{} \Theta$$

3. Эффективность

Определение. Оценка Θ_1^* **не хуже** Θ_2^* , если

$$E(\Theta_1^* - \Theta) \leq E(\Theta_2^* - \Theta)$$

или оценки несмещенные

$$D\Theta_1^* \leq D\Theta_2^*$$

Определение. Оценка Θ^* называется **эффективной** если она не хуже всех остальных оценок.

Теорема 2.1.1. Не существует эффективной оценки в классе всех возможных оценок

Теорема 2.1.2. В классе не смещенных оценок существует эффективная оценка

2.2 Точечные оценки моментов

Определение. Выборочным средним \bar{X}_c называется величина равная среднему арифметическому данных

$$\bar{X}_c = \frac{1}{n} \sum_{i=1}^n X_i$$

Определение. Выборочная дисперсией D_c называется величина

$$D_c = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_c)^2$$

Определение. Исправленной выборочной дисперсией S^2 называется величина

$$S^2 = \frac{n}{n-1} D_c$$

Определение. Выборочным средним квадратическим отклонением σ_c называется величина:

$$\sigma_c = \sqrt{D_c}$$

Определение. Исправленным выборочным средним квадратическим отклонением называется величина:

$$S = \sqrt{S^2}$$

Определение. Выборочным k -ым моментом называется величина

$$\bar{X}^k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

Определение. Модой M_0^* называется величина с наибольшей частотой

Определение. Медианой Me^* вариационного ряда называется значение варианты в середине ряда. Если число четное, то берем среднее от средних.

Теорема 2.2.1. \bar{X}_c является несмещенной состоятельной оценкой для математического ожидания

1. $E\bar{X}_c = EX = a$ — несмещенность
2. $\bar{X}_c \xrightarrow{P} [n \rightarrow \infty] EX = a$

Доказательство. Доделать

□

Теорема 2.2.2. \bar{X}^k является несмещенной состоятельной оценкой для теоретического k -го момента

1. $E\bar{X}^k = EX^k = m_k$ — несмещенность
2. $\bar{X}^k \xrightarrow{P} EX^k = m_k$ — состоятельность

Теорема 2.2.3. D_c и S^2 являются состоятельными оценками для дисперсии

1. D_c — смещенная оценка
2. S^2 — несмещенная оценка

Доказательство. Доделать

□

Примечание. Пусть имеется выборка (X_1, \dots, X_n) неизвестного распределения, знаем тип распределения. Пусть этот тип определяется k неизвестными параметрами. Теоретическое распределение задает теоретический k -тый момент. Например, если распределение непрерывное, то оно задается плотность:

$$f(X, \Theta_1, \dots, \Theta_k) \text{ и } m_k = \int_{-\infty}^{\infty} X^k f(X, \Theta_1, \dots, \Theta_k) dX = f_k(\Theta_1, \dots, \Theta_k)$$

Метод моментов состоит в следующем: вычисляем выборочные моменты и подставляем их в уравнение вместо теоретических. В результате получаем систему уравнений:

$$\begin{cases} \bar{X} = f_1(\Theta_1, \dots, \Theta_k) \\ \bar{X}^2 = f_2(\Theta_1, \dots, \Theta_k) \\ \vdots \\ \bar{X}^k = f_k(\Theta_1, \dots, \Theta_k) \end{cases}$$

Решив эту систему получим оценки параметров $\Theta_1, \dots, \Theta_k$. При этом как правило получаем оценки состоятельные но смещенные

Пример. Пусть случайная величина $X \in U(a, b)$, $a < b$. Обработав стат. данные получили оценки:

$$\bar{X} = 2.25 \quad \bar{X}^2 = 6.75$$

Решение.

$$f(x) = \begin{cases} 0 & , x < a \\ \frac{1}{b-a} & , a \leq x \leq b \\ 0 & , x > b \end{cases}$$

$$EX = \int_a^b X f(X) dX = \dots = \frac{a+b}{2}$$

$$EX^2 = \int_a^b X^2 \frac{1}{b-a} dX = \dots = \frac{a^2 + ab + b^2}{3}$$

$$\begin{cases} 2.25 = \frac{a+b}{2} \\ 6.75 = \frac{a^2+ab+b^2}{3} \end{cases} \Leftrightarrow \dots \Leftrightarrow \begin{cases} b = 4.5 \\ a = 0 \end{cases}$$

Лекция 3

3.1 Метод максимального правдоподобия

Состоит в том, чтобы подобрать параметр таким образом, чтобы вероятность получения данной выборки была наибольшей. Если распределение дискретное, то вероятность выборки

$$P_{\theta}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P_{\theta}(X_1 = x_1)P_{\theta}(X_2 = x_2) \dots P_{\theta}(X_n = x_n)$$

Определение. **Функцией правдоподобия** $L(\bar{X}, \theta)$ называется функция

$P_{\theta}(X_1 = x_1) \dots P_{\theta}(X_n = x_n)$ — в случае дискретного распределения

$$f_{\theta}(x_1) \dots f_{\theta}(x_n) = \prod_{i=1}^n f_{\theta}(x_i) \text{ — в случае абсолютно непрерывного распределения}$$

Определение. **Логарифмической функцией правдоподобия** $M(\bar{X})$ называется функция

$$\ln L(\bar{X})$$

Примечание. Так как логарифм — возрастающая функция, то экстремумы этих функций совпадают.

Определение. **Оценкой максимального правдоподобия** $\hat{\Theta}$ называется значение Θ при котором функция правдоподобия достигает наибольшего значения при фиксированных значениях выборки x_i

Пример. Пусть X_1, \dots, X_n — выборка неизвестного распределения Пуассона с параметром λ .

$$P(X = x_i) = \frac{\lambda^{x_i}}{x_i!} \cdot e^{-\lambda}$$

$$L(\bar{X}) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} \cdot e^{-\lambda} =$$

$$\frac{n \cdot \bar{X}}{\lambda} - n = 0 \implies \lambda = \bar{X}$$

$$\hat{\Theta} = \bar{X}$$

Пример. $f_{a, \sigma^2} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-a)^2}{\sigma^2}}$

$$a = \bar{X} \quad \sigma^2 = D_c$$

Пример. Пусть X_1, \dots, X_n — выборка равномерного распределения вида. $X \in U(0, \Theta)$, $\Theta > 0$. Найти оценки Θ методами моментов и максимального правдоподобия и сравнить их.

$$EX = \frac{a+b}{2} = \frac{\Theta}{2} \implies \hat{\Theta} = 2\bar{X}$$

$$f_{\Theta}(x) = \begin{cases} 0 & x < 0 \\ \frac{1}{\Theta} & 0 \leq x \leq \Theta \\ 0 & x > \Theta \end{cases}$$

$$L(\bar{X}, \Theta) = \prod_{i=1}^n f_{\Theta}(x_i) = \begin{cases} 0 & \Theta < \max(x_i) = X_m \\ \frac{1}{\Theta^n} & \text{иначе} \end{cases}$$

Ясно, что функция $L(\bar{X}, \Theta)$ достигает наибольшего значения при $\Theta = X_m$

Доделать

Примечание. ОМП часто эффективны, но могут быть смещенные.

3.2 Неравенство Раво-Крамера

Пусть известно что случайная величина $X \in F_{\Theta}$ — семейство распределений с параметром Θ

Определение. Носителем семейства распределений F_{Θ} называется множество $C \in \mathbb{R}$ такое что $\forall \Theta : P(X \in C) = 1$

Обозначение.

$$f_{\Theta}(x) = \begin{cases} f_{\Theta}(x) & \text{плотность, если распределение абсолютно непрерывное} \\ P_{\Theta}(X = x) & \text{если распределение дискретное} \end{cases}$$

Определение. Информацией Фишера называется величина

$$I(\Theta) = E\left(\frac{\partial}{\partial \Theta} \ln f_{\Theta}(x)\right)^2$$

при условии если она существует.

Определение. Семейство распределений F_{Θ} называется **регулярным**, если:

1. Существует носитель C семейства F_{Θ} такой, что $\forall x \in C$ функция $\ln f_{\Theta}(x)$ — непрерывно дифференцируема по Θ
2. Информация Фишера $I(\Theta)$ существует и непрерывна по Θ

Теорема 3.2.1 (Неравенство Раво-Крамера). Пусть (X_1, \dots, X_n) — выборка объема n из регулярного семейства распределений F_{Θ} , $\Theta^* = \Theta^*(X_1, \dots, X_n)$ — несмещенная оценка параметра Θ , дисперсия которой ограничена на любом компакте в области Θ . Тогда

$$D\Theta^* \geq \frac{1}{nI(\Theta)}$$

Пример. Пусть $(X_1, \dots, X_n) - \mathbb{X} \in N(a, \sigma^2)$, $a \in \mathbb{R}, \sigma^2 > 0$. Проверим эффективность оценки $a^* = \bar{X}$.

$$f(x) = \frac{\sigma\sqrt{2\pi}}{e} e^{-\frac{(x-a)^2}{\sigma^2}}$$

Возьмем $C = (-\infty, +\infty)$.

$$\ln f(x) = \ln \sigma - \frac{1}{2} \ln -\frac{(x-a)^2}{\sigma^2}$$

$$\frac{\partial}{\partial a} \ln f(x) = \frac{x-a}{\sigma}$$

— непрерывна по a , $\forall a \in \mathbb{R}$.

$$I(a) = E\left(\frac{\partial}{\partial a} \ln f(x)\right)' = \frac{1}{\sigma^2}$$

— непрерывна по a

$$Da^* = D\bar{X} = \frac{1}{n} D\mathbb{X} = \frac{\sigma^2}{n}$$

$$Da^* = \frac{\sigma^2}{n} = \frac{1}{nI(a)} = \frac{1}{n \cdot \frac{1}{\sigma^2}} = \frac{\sigma^2}{n}$$

$$\implies a^* = \bar{X}$$

Примечание. Исправленная дисперсия S^2 также является эффективной оценкой для σ^2 .
BLUE - оценки

$$\Theta^* = \alpha_1 X_1 + \alpha_2 X_2 + \dots + \alpha_n X_n$$

TODO Лекция 4

TODO Лекция 5

TODO Лекция 6

TODO Лекция 7

Лекция 9

8.1 Математическая регрессия

Определение. Регрессией X на Z называется функция $f(z) = E(X|Z = z)$.

$$x = f(z)$$

— уравнение регрессии, а ее график — линия регрессии

Пусть при n экспериментах, при значениях z_1, \dots, z_n случайной величины Z наблюдаемые значения x_1, \dots, x_n случайной величины X .

Обозначим ε разницу между наблюдаемыми и экспериментальными значениями X

$$\varepsilon_i = X_i - E(X|Z = z_i) = X_i - f(z_i)$$

Тогда $x_i = f(z_i) + \varepsilon_i$, $1 \leq i \leq n$, и ε_i — независимые случайные нормальные величины. $\varepsilon_i \in N(0, \sigma^2)$

$$a = 0, \text{ т.к. } E\varepsilon_i = E(X_i) - E(f(z_i)) = EX_i - E(X|Z = z) = 0$$

Цель: по данным значениям как можно более точно оценить функцию $f(z)$

Примечание. При этом предполагается (например из теории), что $f(z)$ — функция определенного типа, параметры которой не известны

8.1.1 Метод наименьшие квадратов

Метод наименьших квадратов состоит в выборе параметров функции $f(z)$, таким образом, чтобы минимизировать сумму квадратов ошибок

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (X_i - f(z_i))^2 \rightarrow \min$$

Определение. Пусть $\Theta = (\Theta_1, \dots, \Theta_n)$ — набор параметров функции $f(z)$.

$$\hat{\Theta} = \min \sum_{i=1}^n \varepsilon_i^2$$

8.2 Линейная парная регрессия

Пусть имеется *линейная регрессия*:

$$f(z) = a + bz$$

Тогда $X_i = a + bz_i + \varepsilon_i$, $1 \leq i \leq n$. Найдем оценки неизвестных параметров a и b методом наименьших квадратов (МНК).

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (X_i - a - bz_i)^2 \rightarrow \min$$

$$\frac{\partial}{\partial a} \sum_{i=1}^n \varepsilon_i^2 = 2 \sum_{i=1}^n (X_i - a - bz_i) \cdot (-1) = -2 \sum_{i=1}^n X_i + 2 \sum_{i=1}^n a + 2b \sum_{i=1}^n z_i = -2(n\bar{X} - na - bn\bar{z}) = 0$$

$$\frac{\partial}{\partial b} \sum_{i=1}^n \varepsilon_i^2 = 2 \sum_{i=1}^n (X_i - a - bz_i) \cdot (-z_i) = -2 \left(\sum_{i=1}^n X_i z_i - a \sum_{i=1}^n z_i - b \sum_{i=1}^n z_i^2 \right) = -2(n\bar{X}\bar{z} - an\bar{z} - bn\bar{z}^2) = 0$$

$$\begin{cases} n\bar{X} - na - bn\bar{z} = 0 \\ n\bar{X}\bar{z} - an\bar{z} - bn\bar{z}^2 \end{cases} \Leftrightarrow \begin{cases} \bar{X} - a - b\bar{z} = 0 \\ \bar{X}\bar{z} - a\bar{z} - b\bar{z}^2 \end{cases} \Leftrightarrow \begin{cases} a + b\bar{z} = \bar{X} \\ a\bar{z} + b\bar{z}^2 = \bar{X}\bar{z} \end{cases}$$

$$\begin{cases} a = \bar{X} - b\bar{z} \\ (\bar{X} - b\bar{z})\bar{z} + b\bar{z}^2 = \bar{X}\bar{z} \end{cases} \Leftrightarrow \begin{cases} a = \bar{x} - b\bar{z} \\ b(\bar{z}^2 - \bar{z}^2) = \bar{X}\bar{z} - \bar{X} \cdot \bar{z} \end{cases} \Leftrightarrow \begin{cases} a = \bar{X} - b\bar{z} \\ b = \frac{\bar{X}\bar{z} - \bar{X} \cdot \bar{z}}{\bar{z}^2 - \bar{z}^2} \end{cases}$$

$$\bar{X}_z = E(\bar{X}|Z = z) = f(z)$$

$$\bar{X}_z = a + bz$$

$$\bar{X}_z = \bar{X} - b\bar{z} + bz$$

$$\bar{X}_z - \bar{X} = b(z - \bar{z})$$

$$\bar{X}_z - \bar{X} = \frac{\bar{X}z - \bar{X} \cdot \bar{z}}{\hat{\sigma}_z^2} \cdot (z - \bar{z})$$

$$\bar{X}_z - \bar{X} = \frac{\bar{X}z - \bar{X} \cdot \bar{z}}{\hat{\sigma}_z \cdot \hat{\sigma}_x} \cdot \frac{\hat{\sigma}_x}{\hat{\sigma}_z} \cdot (z - \bar{z})$$

$$\frac{\bar{X}_z - \bar{X}}{\hat{\sigma}_x} = \frac{\bar{X}z - \bar{X} \cdot \bar{z}}{\hat{\sigma}_z \cdot \hat{\sigma}_x} \cdot \frac{z - \bar{z}}{\hat{\sigma}_z}$$

$$\frac{\bar{X}_z - \bar{X}}{\hat{\sigma}_x} = r_B \cdot \frac{\bar{z} - z}{\hat{\sigma}_z}$$

, где $r_B = \frac{\bar{X}z - \bar{X} \cdot \bar{z}}{\hat{\sigma}_z \cdot \hat{\sigma}_x}$

8.2.1 Геометрический смысл прямой линейной регрессии

Прямая регрессии строится таким образом, чтобы сумма квадратов длин отрезков была наименьшей.

Определение. Выборочный коэффициент линейной корреляции

$$r_B = \frac{\overline{Xz} - \bar{X} \cdot \bar{z}}{\hat{\sigma}_x \cdot \hat{\sigma}_y}$$

Примечание. Отсюда видим, что выборочный коэффициент линейной корреляции является оценкой теоретического коэффициента полученного по методу моментов. Поэтому выборный коэффициент корреляции характеризует силу линейной связи между двумя случайными величинами. Знак r_B показывает является ли корреляция прямой или обратной

8.2.2 Проверка гипотезы по значимости коэффициента линейной корреляции

Пусть двумерная случайная величины распределена нормально. По выборке объема n вычислен r_B , а r — теоретический коэффициент линейной корреляции.

- $H_0 : r = 0$
- против $H_1 : r \neq 0$, т.е. коэффициент r_B статистически значимый

Теорема 8.2.1. Если H_0 верна, то статистика

$$K = \frac{r_B \sqrt{n-2}}{\sqrt{1-r_B^2}} \in T_{n-2}$$

— распределение Стьюдента с $n-2$ степенями свободы

Критерий: Пусть t_{α} — квантиль T_{n-2} уравнения значимости α

8.2.3 Выборочное корреляционное отношение

Пусть имеется k выборок случайной величины X , полученных при уровнях фактора z_1, \dots, z_k . Вычислены общая дисперсия D_O , D_M — межгрупповая дисперсия, D_B — внутригрупповая дисперсия. По теореме о разложении дисперсии $D_O = D_M + D_B$ Доделать