

MS L_EC 04

isagila

ablearthy

Собрано 07.08.2024 в 15:41



Содержание

1. Лекции	3
1.1. Лекция 24.02.08.	3
1.2. Лекция 24.02.15.	4
1.3. Лекция 24.02.22.	7
1.4. Лекция 24.02.29.	10
1.5. Лекция 24.03.07.	14
1.6. Лекция 24.03.14.	17
1.7. Лекция 24.03.21.	20
1.8. Лекция 24.03.28.	23
1.9. Лекция 24.04.04.	26
1.10. Лекция 24.04.11.	29
1.11. Лекция 24.04.18.	33
1.12. Лекция 24.04.25.	37
1.13. Лекция 24.05.02.	40
1.14. Лекция 24.05.16.	43
1.15. Лекция 24.05.23.	45
1.16. Лекция 24.05.30.	48

1. Лекции

1.1. Лекция 24.02.08.

Выборки

Def 1.1.1. Под генеральной совокупностью понимают все результаты данной серии экспериментов (или экспериментальных значений случайной величины).

Def 1.1.2. Под выборочной совокупностью понимают имеющиеся у нас данные (выборка из генеральной совокупности, возможно неполная).

Замечание 1.1.3. Выборочная совокупность не всегда отражает реальное поведение случайной величины. В качестве иллюстрации можно рассматривать известный пример с бронированием самолетов (классическая ошибка выжившего).

Def 1.1.4. Репрезентативной выборкой называется выборка, имеющая то же самое распределение, что и у генеральной совокупности.

Замечание 1.1.5. В дальнейшем в курсе предполагаем, что все выборки репрезентативные.

Def 1.1.6 (Первое). Выборкой объема n называется набор экспериментальных данных (x_1, \dots, x_n) .

Def 1.1.7 (Второе). Выборкой объема n называется набор X_1, \dots, X_n независимых одинаково распределенных случайных величин.

Замечание 1.1.8. Таким образом $\mathbb{E}(X_i) = \mathbb{E}(X_1)$ и $\mathbb{D}(X_i) = \mathbb{D}(X_1)$, поэтому обычно числовые характеристики i -того заменяются числовыми характеристиками первого экземпляра.

Выборочные характеристики

Выборку можно рассматривать как дискретную случайную величину.

X	x_1	\dots	x_n
p^*	$\frac{1}{n}$	\dots	$\frac{1}{n}$

Def 1.1.9. Средним выборочным \bar{x} называется число

$$\bar{x} = \frac{x_1 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

Это будет оценкой неизвестного математического ожидания.

Def 1.1.10. Выборочной дисперсией \mathbb{D}^* называется число

$$\mathbb{D}^* = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Это будет оценкой неизвестной дисперсии.

Def 1.1.11. Выборочным среднеквадратическим отклонением называется число

$$\sigma^* = \sqrt{\mathbb{D}^*}$$

Это будет оценкой неизвестного среднеквадратического отклонения.

Def 1.1.12. Выборочной функцией распределения $F^*(y)$ называется функция

$$F^*(y) = \frac{\text{число данных } x_i \in (-\infty; y)}{n}$$

Теорема 1.1.13. (Гливленко—Кантелли) Пусть имеется выборка $\vec{X} = (x_1, \dots, x_n)$ объема n . Обозначим $F^*(y)$ — эмпирическую функцию распределения, а $F(y)$ — теоретическую функцию распределения. Тогда

$$\sup_{y \in \mathbb{R}} |F^*(y) - F(y)| \xrightarrow{P} 0 \quad \text{при } n \rightarrow \infty$$

Начальная обработка статданных

I. Ранжирование выборки

Упорядочиваем данные по возрастанию, в результате получаем вариационный ряд вида $X_{(1)}, \dots, X_{(n)}$.

Def 1.1.14. Разность $X_{(n)} - X_{(1)}$ называется размахом выборки.

Def 1.1.15. Элемент $X_{(i)}$ в полученном ряде называется i -той порядковой статистикой.

Замечание 1.1.16. Если мы при этом объединяем повторяющиеся результаты (с учетом числа этих результатов), то получаем частотный вариационный ряд.

II. Разбиение на интервалы

Если много неповторяющихся данных, то разбиваем их на интервалы и составляем интервальный вариационный ряд.

Замечание 1.1.17. Есть два подхода к разбиению на интервалы.

1. Берем интервалы одинаковой длины. Это удобнее для построения гистограммы и выдвижения гипотезы и типа распределения.
2. Берем равнонаполненные интервалы. Это удобнее при проверке гипотез о типе распределения.

Замечание 1.1.18. Обычно (но далеко не всегда) количество интервалов вычисляется по формуле Стерджеса.

$$K \approx 1 + \log_2 n$$

В результате получаем K интервалов $[a_{i-1}; a_i)$ и считаем их частоты v_i — число данных, попавших в i -ый интервал. Величина $\frac{v_i}{n}$ называется относительной частотой (она является оценкой теоретической вероятности попадания случайной величины в данный интервал). Если заменяем интервалы на их середины $c_i = \frac{a_{i-1} + a_i}{2}$, то опять получаем огрубленный вариационный ряд, по которому можно дать точные оценки числовым характеристикам распределения. Для этого можно использовать формулы

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n c_i v_i \quad \mathbb{D}^* = \frac{1}{n} \sum_{i=1}^n (c_i - \bar{x})^2 v_i$$

Геометрическая интерпретация данных

Обычно удобнее визуализировать данные в виде гистограммы. На плоскости строим набор прямоугольников для каждого интервала. Основание это $[a_{i-1}; a_i)$ длины $l = a_i - a_{i-1}$, а высоту берем пропорционально частоте, причем таким образом, чтобы общая площадь фигуры равнялась единице, т.е. высота будет равна $\frac{v_i}{nl}$. Гистограмма является приближением плотности распределения, и по ее виду можно выдвинуть гипотезы о типе распределения.

Теорема 1.1.19. Если число интервалов $K(n) \rightarrow \infty$ и $\frac{K(n)}{n} \rightarrow 0$ при $n \rightarrow \infty$, то гистограмма по вероятности поточечно сходится к теоретической плотности.

1.2. Лекция 24.02.15.

Точечные оценки

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ объема n — набор независимых экземпляров случайной величины X .

Def 1.2.1. Статистикой называется измеримая функция $\Theta^* = \Theta^*(X_1, \dots, X_n)$.

Пусть требуется найти приближенную оценку неизвестного параметра Θ по выборке (X_1, \dots, X_n) . Оценка считается при помощи некоторой статистики $\Theta^* = \Theta^*(X_1, \dots, X_n)$.

Свойства статистических оценок

Def 1.2.2. Статистика $\Theta^* = \Theta^*(X_1, \dots, X_n)$ неизвестного параметра Θ называется состоятельной, если $\Theta^* \xrightarrow{P} \Theta$ при $n \rightarrow \infty$.

Def 1.2.3. Статистика $\Theta^* = \Theta^*(X_1, \dots, X_n)$ неизвестного параметра Θ называется несмещенной, если $\mathbb{E}(\Theta^*) = \Theta$.

Замечание 1.2.4. Оценка называется асимптотически несмещенной, если $\mathbb{E}(\Theta^*) \rightarrow \Theta$ при $n \rightarrow \infty$.

Def 1.2.5. Оценка Θ_1^* не хуже оценки Θ_2^* , если

$$\mathbb{E}((\Theta_1^* - \Theta)^2) \leq \mathbb{E}((\Theta_2^* - \Theta)^2)$$

Если Θ_1^* и Θ_2^* — несмещенные оценки, то это равносильно тому, что $\mathbb{D}(\Theta_1^*) \leq \mathbb{D}(\Theta_2^*)$.

Def 1.2.6. Оценка Θ^* называется эффективной, если она не хуже всех остальных оценок.

Замечание 1.2.7. В классе всех возможных оценок не существует эффективной оценки.

Теорема 1.2.8. В классе несмещенных оценок существует эффективная оценка, причем единственная.

Точечные оценки моментов

Def 1.2.9. Выборочным средним \bar{x} называется величина

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Def 1.2.10. Выборочной дисперсией \mathbb{D}^* называется величина

$$\mathbb{D}^* = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Def 1.2.11. Исправленной выборочной дисперсией S^2 называется величина

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n}{n-1} \mathbb{D}^*$$

Def 1.2.12. Выборочным среднеквадратичным отклонением называется величина

$$\sigma^* = \sqrt{\mathbb{D}^*}$$

Def 1.2.13. Исправленным выборочным среднеквадратичным отклонением называется величина

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

Def 1.2.14. Выборочным k -ым моментом $\overline{x^k}$ называется величина

$$\overline{x^k} = \frac{1}{n} \sum_{i=1}^n x_i^k$$

Def 1.2.15. Выборочной модой Mo^* называется варианта с наибольшей частотой.

$$Mo^* = x_i \mid n_i = \max(n_1, \dots, n_k)$$

Def 1.2.16. Выборочной медианой Me^* называется значение варианты x_i в середине ряда.

$$Me^* = \begin{cases} x_{(k)}, & \text{если } n = 2k - 1 \\ \frac{1}{2} (x_{(k)} + x_{(k+1)}), & \text{если } n = 2k \end{cases}$$

Теорема 1.2.17. Выборочное среднее \bar{x} является несмещенной состоятельной оценкой для математического ожидания.

$$\mathbb{E}(\bar{x}) = \mathbb{E}(X) = a \quad \bar{x} \xrightarrow{P} \mathbb{E}(X) = a$$

□ Покажем несмещенность.

$$\mathbb{E}(\bar{x}) = \mathbb{E}\left(\frac{x_1 + \dots + x_n}{n}\right) = \frac{1}{n} \cdot n \cdot \mathbb{E}(x_1) = \mathbb{E}(x_1)$$

Покажем состоятельность.

$$\bar{x} = \frac{x_1 + \dots + x_n}{n} \xrightarrow{P} \mathbb{E}(X) \text{ при } n \rightarrow \infty$$

Это закон больших чисел. ■

Теорема 1.2.18. Выборочный k -ый момент $\overline{x^k}$ является несмещенной состоятельной оценкой для теоретического k -ого момента.

$$\mathbb{E}(\overline{x^k}) = \mathbb{E}(X^k) = m_k \quad \overline{x^k} \xrightarrow{P} \mathbb{E}(X^k) = m_k \text{ при } n \rightarrow \infty$$

□ Это следует из предыдущей теоремы, если в качестве случайной величины взять x^k . ■

Теорема 1.2.19. Выборочные дисперсии \mathbb{D}^* и S^2 являются состоятельными оценками для дисперсии. При этом \mathbb{D}^* — смещенная оценка (есть систематическая вниз), а S^2 — несмещенная оценка.

□ Заметим, что

$$\mathbb{D}^* = \frac{1}{n} \sum (x_i - \bar{x})^2 = \overline{x^2} - (\bar{x})^2$$

$$\mathbb{D}(\bar{x}) = \mathbb{E}((\bar{x})^2) - (\mathbb{E}(\bar{x}))^2 \implies \mathbb{E}((\bar{x})^2) = (\mathbb{E}(\bar{x}))^2 + \mathbb{D}(\bar{x})$$

Таким образом

$$\mathbb{E}(\mathbb{D}^*) = \mathbb{E}(\overline{x^2}) - \mathbb{E}((\bar{x})^2) = \mathbb{E}(x^2) - ((\mathbb{E}(x))^2 + \mathbb{D}(\bar{x})) = (\mathbb{E}(x^2) - (\mathbb{E}(x))^2) - \mathbb{D}(\bar{x}) = \mathbb{D}(x) - \mathbb{D}(\bar{x})$$

Видно, что будет смещение на величину $\mathbb{D}(\bar{x})$. Преобразуем полученное выражение в более удобную форму.

$$\mathbb{D}(x) - \mathbb{D}(\bar{x}) = \mathbb{D}(x) - \frac{1}{n^2} \sum_{i=1}^n \mathbb{D}(x_i) = \mathbb{D}(x) - \frac{1}{n} \cdot \mathbb{D}(x) = \frac{n-1}{n} \cdot \mathbb{D}(x)$$

Значит

$$\mathbb{E}(S^2) = \mathbb{E}\left(\frac{n}{n-1} \mathbb{D}^*\right) = \frac{n}{n-1} \cdot \frac{n-1}{n} \cdot \mathbb{D}(x) = \mathbb{D}(x)$$

Т.е. S^2 это несмещенная оценка. Далее покажем состоятельность оценок.

$$\mathbb{D}^* = \overline{x^2} - (\bar{x})^2 \xrightarrow{P} \mathbb{E}(x^2) - (\mathbb{E}(x))^2 = \mathbb{D}(x)$$

$$S^2 = \underbrace{\frac{n}{n-1}}_{\rightarrow 1} \mathbb{D}^* \xrightarrow{P} \mathbb{D}(x)$$

■

Замечание 1.2.20. Т.к. $\frac{n}{n-1} \rightarrow 1$ при $n \rightarrow \infty$, то $\mathbb{E}(\mathbb{D}^*) \xrightarrow{P} \mathbb{D}(x)$ при $n \rightarrow \infty$. Значит выборочная дисперсия является асимптотически несмещенной оценкой, поэтому на практике при $n \geq 100$ можно считать обычную выборочную дисперсию, а при $n < 100$ следует ее заменить на исправленную выборочную дисперсию.

Метод моментов (Пирсона)

Зная выборочные моменты можно дать оценки остальным параметрам распределения. Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ неизвестного распределения, но при этом знаем, что данное распределение определенного типа, задаваемого k параметрами $\Theta = (\Theta_1, \dots, \Theta_k)$. Зная параметры можем вычислить теоретические k -ые моменты. Например, если распределение непрерывное, то

$$m_i = \int_{\mathbb{R}} x^k f(x, \Theta_1, \dots, \Theta_k) dx = h_i(\Theta_1, \dots, \Theta_k)$$

Вычислим выборочные моменты и подставим в эти формулы, получим систему уравнений.

$$\begin{cases} \bar{x} = h_1(\Theta_1, \dots, \Theta_k) \\ \overline{x^2} = h_2(\Theta_1, \dots, \Theta_k) \\ \vdots \\ \overline{x^k} = h_k(\Theta_1, \dots, \Theta_k) \end{cases}$$

Решив эту систему, находим оценки $\Theta_1^*, \dots, \Theta_k^*$ неизвестных параметров.

Замечание 1.2.21. При этом обычно получаем состоятельные оценки, но смещенные.

Пример 1.2.22. Пусть $X \in U(a; b)$. При обработке статданных получили оценки первого и второго момента $\bar{x} = 2.25$ и $\overline{x^2} = 6.75$. Необходимо дать оценки неизвестных параметров a и b .

Решение: Плотность будет иметь вид

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

Тогда

$$\mathbb{E}(x) = \frac{a+b}{2} \quad \mathbb{E}(x^2) = \frac{a^2 + ab + b^2}{3}$$

Составим систему

$$\begin{cases} 2.25 = \frac{a^* + b^*}{2} \\ 6.75 = \frac{(a^*)^2 + a^*b^* + (b^*)^2}{3} \end{cases} \iff \begin{cases} a^* + b^* = 4.5 \\ (a^*)^2 + a^*b^* + (b^*)^2 = 20.25 \end{cases} \iff \begin{cases} a^* + b^* = 4.5 \\ a^*b^* = 0 \end{cases} \xLeftrightarrow{a \leq b} \begin{cases} a^* = 0 \\ b^* = 4.5 \end{cases}$$

1.3. Лекция 24.02.22.

Пусть имеется выборка распределения известного нам типа, но с неизвестными параметрами. Требуется найти эти параметры. На прошлой лекции был рассмотрен метод моментов. Его недостаток заключается в том, что он не гарантирует эффективность оценок.

Метод максимального правдоподобия (Фишера)

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ распределения известного типа, задаваемого параметрами $\Theta = (\Theta_1, \dots, \Theta_k)$. Идея метода заключается в подборе параметров таким образом, чтобы вероятность получения данной выборки при случайном эксперименте была наибольшей. Например, если распределение дискретное, то эта вероятность будет равна

$$P_{\Theta}(x_1, \dots, x_n) = P_{\Theta}(x = x_1) \cdot \dots \cdot P_{\Theta}(x = x_n)$$

Def 1.3.1. Функцией правдоподобия $L(\vec{X}, \Theta)$ называется функция

$$L(\vec{X}, \Theta) = \prod_{i=1}^n P(x = x_i) \quad \text{при дискретном распределении}$$

$$L(\vec{X}, \Theta) = \prod_{i=1}^n f_{\Theta}(x_i) \quad \text{при непрерывном распределении}$$

Def 1.3.2. Логарифмической функцией правдоподобия называется функция $\ln L(\vec{X}, \Theta)$.

Замечание 1.3.3. Т.к. $y = \ln x$ возрастающая функция, то их точки экстремума совпадают, но искать их проще во втором случае.

Def 1.3.4. Оценкой максимального правдоподобия $\hat{\Theta}$ называется значение Θ , при котором функция правдоподобия достигает наибольшего значения (при фиксированных значениях x_1, \dots, x_n).

Пример 1.3.5. Пусть $\vec{X} = (X_1, \dots, X_n)$ из распределения Пуассона с неизвестным параметром $\lambda > 0$. Требуется найти оценку максимального правдоподобия параметра λ .

Решение: Составим функцию правдоподобия.

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

Перейдем к логарифмической функции правдоподобия.

$$\ln L(\vec{X}, \lambda) = n\bar{x} \ln \lambda - \ln \prod_{i=1}^n x_i! - n\lambda$$

Вычислим производную по λ и приравняем ее к нулю, чтобы найти экстремум.

$$\frac{d}{d\lambda} (\ln L(\vec{X}, \lambda)) = \frac{n\bar{x}}{\lambda} - n = 0 \implies \lambda = \bar{x}$$

Покажем, что это действительно точка максимума.

$$\frac{d^2}{d\lambda^2} (\ln L(\vec{X}, \lambda)) = -\frac{n\bar{x}}{\lambda^2} < 0$$

Итого $\lambda = \bar{x}$ это оценка максимального правдоподобия.

Пример 1.3.6. Пусть $\vec{X} = (X_1, \dots, X_n)$ — выборка из $N(a; \sigma^2)$, где $a \in \mathbb{R}$ и $\sigma > 0$. Требуется найти оценку максимального правдоподобия.

Решение: Составим функцию правдоподобия.

$$L(\vec{X}, a, \sigma) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x_i - a)^2}{2\sigma^2}\right) = \frac{1}{\sigma^n (2\pi)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - a)^2\right)$$

Перейдем к логарифмической функции правдоподобия.

$$\ln L(\vec{X}, a, \sigma) = -n \ln \sigma - \frac{n}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - a)^2$$

Найдем частные производные по a и σ .

$$\begin{aligned}\frac{\partial}{\partial a} \ln L(\vec{X}, a, \sigma) &= -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - a) \cdot (-1) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - a) = \frac{n}{\sigma^2} (\bar{x} - a) \\ \frac{\partial}{\partial \sigma} \ln L(\vec{X}, a, \sigma) &= -\frac{n}{\sigma} - \frac{1}{2} \cdot (-2) \sigma^{-3} \cdot \sum_{i=1}^n (x_i - a)^2 = \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - a)^2 - \frac{n}{\sigma}\end{aligned}$$

Приравняем их к нулю. Из производной по a получаем $\hat{a} = \bar{x}$. Поработаем с производной по σ .

$$\frac{1}{\sigma^3} \sum_{i=1}^n (x_i - a)^2 - \frac{n}{\sigma} = 0 \implies \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - a)^2 \implies \hat{\sigma}^2 = \mathbb{D}^*$$

То, что это действительно точка максимума можно показать с помощью вторых дифференциалов (самостоятельно).

Пример 1.3.7. Пусть $\vec{X} = (X_1, \dots, X_n)$ из $U(0; \Theta)$, где $\Theta > 0$. Найти оценки Θ методом моментов и методом максимального правдоподобия. Сравнить результаты.

Решение: По методу моментов

$$\mathbb{E}(X) = \frac{0 + \Theta}{2} \implies \bar{x} = \frac{\Theta^*}{2} \implies \Theta^* = 2\bar{x}$$

Далее рассмотрим метод максимального правдоподобия. Обозначим $X_{(n)}$ n -ую порядковую статистику (т.е. максимальный элемент выборки), тогда

$$L(\vec{X}, \Theta) = \prod_{i=1}^n f_{\Theta}(x_i) = \begin{cases} 0, & \Theta < X_{(n)} \\ \frac{1}{\Theta^n}, & \Theta \geq X_{(n)} \end{cases}$$

Значит $\hat{\Theta} = X_{(n)}$ это оценка максимального правдоподобия.

Видим, что полученные оценки отличаются. Сравним их. Сначала рассмотрим оценку по методу моментов. Отметим, что она будет несмещенной, т.к.

$$\Theta^* = 2\bar{x} \quad \mathbb{E}(\Theta^*) = 2\mathbb{E}(\bar{x}) = 2\mathbb{E}(x) = \Theta$$

Таким образом

$$\mathbb{E}((\Theta^* - \Theta)^2) = \mathbb{D}(\Theta^*) = 4\mathbb{D}(\bar{x}) = 4\frac{\mathbb{D}(x)}{n} = \frac{4}{n} \cdot \frac{\Theta^2}{12} = \frac{\Theta^2}{3n}$$

Далее рассмотрим оценку по методу максимального правдоподобия. Вначале рассмотрим случайную величину $X_{(n)} = \max(X_1, \dots, X_n)$. Найдем ее функцию распределения.

$$F_{X_{(n)}} = \mathbb{P}(X_{(n)} < x) = \mathbb{P}(X_1 < x, \dots, X_n < x) = \mathbb{P}(X_1 < x) \cdot \dots \cdot \mathbb{P}(X_n < x) = (F(x))^n = \begin{cases} 0, & x < 0 \\ \frac{x^n}{\Theta^n}, & 0 \leq x \leq \Theta \\ 1, & x > \Theta \end{cases}$$

Далее найдем плотность.

$$f_{X_{(n)}}(x) = F'_{X_{(n)}}(x) = \begin{cases} 0, & x < 0 \\ \frac{nx^{n-1}}{\Theta^n}, & 0 \leq x \leq \Theta \\ 0, & x > \Theta \end{cases}$$

Вычислим первый момент.

$$\mathbb{E}(X_{(n)}) = \int_0^{\Theta} x \cdot \frac{nx^{n-1}}{\Theta^n} dx = \frac{1}{\Theta^n} \int_0^{\Theta} nx^n dx = \frac{1}{\Theta^n} \frac{nx^{n+1}}{n+1} \Big|_0^{\Theta} = \frac{1}{\Theta^n} \cdot \frac{n\Theta^{n+1}}{n+1} = \frac{n}{n+1} \Theta$$

Значит эта оценка будет смещенной вниз. Подправим ее следующим образом

$$\tilde{\Theta} = \frac{n+1}{n} \hat{\Theta} = \frac{n+1}{n} X_{(n)}$$

Такая оценка будет уже несмещенной. Вычислим ее второй момент.

$$\mathbb{E}(\tilde{\Theta}^2) = \frac{(n+1)^2}{n^2} \int_0^{\Theta} x^2 \frac{nx^{n-1}}{\Theta^n} dx = \frac{(n+1)^2}{n} \cdot \frac{1}{\Theta^n} \int_0^{\Theta} x^{n+1} dx = \frac{(n+1)^2}{n\Theta^n} \cdot \frac{x^{n+2}}{n+2} \Big|_0^{\Theta} = \frac{(n+1)^2 \Theta^2}{n(n+2)}$$

Теперь вычислим дисперсию этой несмещенной оценки.

$$\mathbb{E}((\tilde{\Theta} - \Theta)^2) = \mathbb{D}(\tilde{\Theta}) = \mathbb{E}(\tilde{\Theta}^2) - (\mathbb{E}(\tilde{\Theta}))^2 = \frac{(n+1)^2 \Theta^2}{n(n+2)} - \Theta^2 = \frac{n^2 + 2n + 1 - n^2 - 2n}{n(n+2)} \cdot \Theta^2 = \frac{\Theta^2}{n(n+2)}$$

Итак, теперь можно сравнить эти две оценки.

$$\mathbb{D}(\Theta^*) = \frac{\Theta^2}{3n} \quad \mathbb{D}(\tilde{\Theta}) = \frac{\Theta^2}{n(n+2)}$$

Хорошо видно, что вторая оценка сходится быстрее, т.е. эффективность оценки по методу максимального правдоподобия заметно выше. Оценка по методу моментов сходится со скоростью $\sim \frac{1}{\sqrt{n}}$, а по методу максимального правдоподобия со скоростью $\sim \frac{1}{n}$.

Замечание 1.3.8. Т.к. $\mathbb{E}(X) = \frac{\Theta}{2}$, то из примера выше следует, что оценка \bar{x} не является эффективной оценкой математического ожидания. В данном случае оценка $\widetilde{\mathbb{E}(X)} = \frac{n+1}{2n} X_{(n)}$ будет лучше. Можно показать, что эта оценка является эффективной.

Замечание 1.3.9. В общем случае равномерного распределения несмещенной эффективной оценкой математического ожидания будет $\frac{1}{2}(X_{(1)} + X_{(n)})$. А несмещенной эффективной оценкой длины интервала будет $\frac{n+1}{n-1}(X_{(n)} - X_{(1)})$.

Замечание 1.3.10. Таким образом мы видим, что при равномерном распределении нельзя откидывать крайние значения выборки.

Замечание 1.3.11. При методе максимального правдоподобия как правило получаем состоятельные и эффективные оценки.

Неравенство Рао—Крамера

Пусть известно, что случайная величина $X \in \mathcal{F}_\Theta$, где \mathcal{F}_Θ это семейство распределений с параметром Θ .

Def 1.3.12. Носителем семейства распределений \mathcal{F}_Θ называется множество $C \in \mathbb{R}$ такое, что $\forall \Theta \mid \mathbb{P}(X \in C) = 1$.

Замечание 1.3.13. Далее для того, чтобы не формулировать теорему дважды, под плотностью будем иметь в виду

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

Def 1.3.14. Информацией Фишера $I(\Theta)$ семейства распределений \mathcal{F}_Θ называется величина

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

при условии, что она существует.

Def 1.3.15. Семейство распределений \mathcal{F}_Θ называется регулярным, если

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

Теорема 1.3.16. (Неравенство Рао—Крамера) Пусть (X_1, \dots, X_n) это выборка из регулярного семейства \mathcal{F}_Θ , $\Theta^* = \Theta^*(X_1, \dots, X_n)$ — несмещенная оценка параметра Θ , дисперсия $\mathbb{D}(\Theta^*)$ ограничена на любом компакте (компакт \approx ограниченное замкнутое множество). Тогда

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

Замечание 1.3.17. Если $\mathbb{D}(\Theta^*) = \frac{1}{nI(\Theta)}$, то оценка Θ^* будет эффективной.

Пример 1.3.18. Пусть $\vec{X} = (X_1, \dots, X_n)$ — выборка из $N(a; \sigma^2)$, где $a \in \mathbb{R}$ и $\sigma > 0$. Проверим эффективность оценки $a^* = \bar{x}$.

Решение: В качестве носителя выберем $C = (-\infty; \infty)$. Далее вычислим $\ln f(x)$.

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-a)^2}{2\sigma^2}\right) \Rightarrow \ln f(x) = -\ln \sigma - \ln \sqrt{2\pi} - \frac{(x-a)^2}{2\sigma^2}$$

Теперь вычислим частную производную по a .

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

Видим, что эта производная непрерывна по a . Значит первое условие из определения 1.3.15 выполнено. Теперь вычислим информацию Фишера.

$$I(\Theta) = \mathbb{E} \left(\left(\frac{\partial}{\partial a} \ln f(x) \right)^2 \right) = \mathbb{E} \left(\left(\frac{X-a}{\sigma^2} \right)^2 \right) = \frac{\mathbb{E}((X - \mathbb{E}(X))^2)}{\sigma^4} = \frac{\mathbb{D}(X)}{\sigma^4} = \frac{1}{\sigma^2}$$

Видим, что информация Фишера существует и непрерывна по a . Значит выполнено второе условие из определения 1.3.15. Таким образом семейство нормальных распределений $N(a; \sigma^2)$ регулярно относительно параметра a . Аналогично можно показать, что оно будет регулярно относительно параметра σ . Вычислим дисперсию оценки.

$$\mathbb{D}(a^*) = \mathbb{D}(\bar{x}) = \frac{\mathbb{D}(X)}{n} = \frac{\sigma^2}{n}$$

В правой части неравенства 1.3.16 имеем

$$\frac{1}{nI(\Theta)} = \frac{\sigma^2}{n}$$

Видим, что дисперсия оценки получилась минимально возможной. Значит \bar{x} это эффективная оценка параметра a . Аналогично можно показать, что S^2 — несмещенная эффективная оценка σ^2 .

Замечание 1.3.19. В практической статистике часто используют линейные оценки $\Theta^* = \sum c_i x_i$. Лучшая оценка в таком классе называется наилучшей линейной несмещенной оценкой или BLUE-оценкой. В примере с нормальным распределением оценка $a^* = \bar{x}$ будет BLUE-оценкой, а в примере с равномерным распределением оценка $\Theta^* = 2\bar{x}$ будет BLUE-оценкой.

1.4. Лекция 24.02.29.

Основные распределения математической статистики

Def 1.4.1. Случайная величина имеет нормальное распределение $N(a; \sigma^2)$ с параметрами $a \in \mathbb{R}$ и $\sigma^2 (\sigma > 0)$, если его плотность имеет вид

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-a)^2}{2\sigma^2}\right) \quad x \in \mathbb{R}$$

Def 1.4.2. Распределение $N(0; 1)$ с параметрами $a = 0, \sigma = 1$ называется стандартным нормальным распределением. Его плотность имеет вид

$$f(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

Свойства нормального распределения

1. Смысл параметров

$$\mathbb{E}(X) = a \quad \mathbb{D}(X) = \sigma^2$$

2. Линейность

$$\xi \in N(a; \sigma^2) \implies b\xi + c \in N(ab + c; b^2\sigma^2)$$

3. Стандартизация

$$\xi \in N(a; \sigma^2) \implies \frac{\xi - a}{\sigma} \in N(0; 1)$$

4. Устойчивость относительно суммирования

$$\left. \begin{array}{l} \xi_1 \in N(a_1; \sigma_1^2) \\ \xi_2 \in N(a_2; \sigma_2^2) \end{array} \right\} \implies \xi_1 + \xi_2 \in N(a_1 + a_2; \sigma_1^2 + \sigma_2^2)$$

I. Распределение «хи—квадрат»

Def 1.4.3. Распределением «хи—квадрат» H_k с k степенями свободы называется распределение суммы k квадратов независимых стандартных нормальных величин.

$$\chi_k^2 = X_1^2 + \dots + X_k^2 \quad X_i \in N(0; 1) \text{ независимые}$$

Свойства распределения «хи—квадрат»

Lm 1.4.4.

$$\mathbb{E}(\chi_k^2) = k$$

□

$$\begin{aligned}\mathbb{D}(X_1) &= \mathbb{E}(X_1^2) - (\mathbb{E}(X_1))^2 = \mathbb{E}(X_1^2) = 1 \\ \mathbb{E}(\chi_k^2) &= \mathbb{E}(X_1^2 + \dots + X_k^2) = k\mathbb{E}(X_1^2) = k\end{aligned}$$

■

Замечание 1.4.5. Распределение «хи-квадрат» устойчиво относительно суммирования.

$$\left. \begin{array}{l} X_1 \in \chi_n^2 \\ X_2 \in \chi_m^2 \\ X_1, X_2 \text{ независимы} \end{array} \right\} \Rightarrow X_1 + X_2 \in \chi_{n+m}^2$$

II. Распределение Стюдента

Def 1.4.6. Пусть случайные величины X_0, \dots, X_k независимы и имеют стандартное нормальное распределение. Распределением Стюдента T_k с k степенями свободы называется распределение случайной величины

$$t_k = \frac{X_0}{\sqrt{\frac{1}{k}(X_1^2 + \dots + X_k^2)}} = \frac{X_0}{\sqrt{\frac{1}{k}\chi_k^2}}$$

Свойства распределения Стюдента

Замечание 1.4.7.

$$\mathbb{E}(t_k) = 0$$

Lm 1.4.8.

$$t_k \Rightarrow N(0; 1)$$

□

$$\frac{\chi_k^2}{k} \xrightarrow{P} \mathbb{E}(X_1^2) = 1 \Rightarrow t_k = \frac{X_0}{\sqrt{\frac{1}{k}\chi_k^2}} \xrightarrow{P} X_0$$

■

III. Распределение Фишера—Снедекора (F—распределение)

Def 1.4.9. Распределение Фишера—Снедекора $F_{m,n}$ с m и n степенями свободы называется случайной величины

$$f_{m,n} = \frac{n\chi_m^2}{m\chi_n^2}$$

Замечание 1.4.10.

$$\mathbb{E}(f_{m,n}) = \frac{n}{n-2}$$

Преобразование нормальных выборок

Пусть $\vec{X} = (X_1, \dots, X_n)^T$ — выборка из $N(0; 1)$, т.е. $X_i \in N(0; 1)$ — независимы. A это невырожденная матрица порядка n . Рассмотрим $\vec{Y} = A\vec{X}$, где $Y_i = a_{i,1}X_1 + \dots + a_{i,n}X_n$. По свойствам нормального распределения эти компоненты будут нормальными случайными величинами, но в общем случае зависимыми. Нас в основном интересует случай, когда A — ортогональная матрица.

Многомерное нормальное распределение

Def 1.4.11. Пусть случайный вектор $\vec{\xi} = (\xi_1, \dots, \xi_n)$ имеет вектор средних $\mathbb{E}(\vec{\xi}) = \vec{a} = (\mathbb{E}(\xi_1), \dots, \mathbb{E}(\xi_n))^T$, а K — симметричная положительно определенная матрица. Вектор $\vec{\xi}$ имеет многомерное нормальное распределение в \mathbb{R}^n с параметрами \vec{a} и K , если его плотность имеет вид

$$f_{\vec{a}, K}(\vec{x}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{\det K}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{a})^T K^{-1}(\vec{x} - \vec{a})\right)$$

Многомерное нормальное распределение обозначается как $\vec{\xi} \in N(\vec{a}; K)$.

Замечание 1.4.12. Матрица $K = \mathbb{D}(\vec{\xi}) = (\text{cov}(\xi_i, \xi_j))$ — матрица ковариаций.

Lm 1.4.13. Если $K = E$ и $\vec{a} = \vec{0}$, то имеем вектор из независимых стандартных нормальных случайных величин.

□ Запишем плотность с данными параметрами

$$f_{\xi}(\vec{x}) = \frac{1}{(\sqrt{2\pi})^n} \exp\left(-\frac{1}{2} \begin{pmatrix} \vec{x} \\ x \end{pmatrix}^T \begin{pmatrix} \vec{x} \\ x \end{pmatrix}\right) = \frac{1}{(\sqrt{2\pi})^n} \exp\left(-\frac{1}{2} (x_1^2 + \dots + x_n^2)\right) = f_{\xi_1}(x_1) \cdot \dots \cdot f_{\xi_n}(x_n)$$

где $f_{\xi_i}(x_i)$ — плотность стандартного нормального распределения. Заметим, что случайные величины ξ_i независимы, т.к. плотность совместного распределения равна произведению частных плотностей. ■

Замечание 1.4.14. Пусть \vec{X} состоит из независимых нормальных стандартных случайных величин, B — невырожденная матрица порядка n . Тогда $\vec{Y} = B\vec{X} + \vec{a}$ имеет многомерное нормальное распределение с параметрами \vec{a} и $K = BB^T$. Это эквивалентное определение многомерных нормальных распределений — все они получаются таким образом.

Замечание 1.4.15. Случайный вектор имеет многомерное нормальное распределение, если всего его компоненты это нормальные случайные величины и нет функциональной зависимости одной компоненты от остальных.

Лм 1.4.16. Пусть случайный вектор \vec{X} состоит из независимых стандартных нормальных случайных величин, а C это ортогональная матрица. Тогда $\vec{Y} = C\vec{X}$ также состоит из независимых стандартных нормальных случайных величин.

□ Т.к. C — ортогональная матрица, то $C^{-1} = C^T$, поэтому $K = CC^{-1} = E$. По лемме 1.4.13 получаем, что \vec{Y} состоит из независимых стандартных нормальных случайных величин. ■

Лм 1.4.17. Пусть случайный вектор $\vec{\xi}$ имеет многомерное нормальное распределение с параметрами \vec{a} и K . Тогда $\vec{\eta} = B^{-1}(\vec{\xi} - \vec{a})$, где $B = \sqrt{K}$, состоит из независимых стандартных нормальных случайных величин.

Лм 1.4.18. Пусть случайный вектор $\vec{\xi}$ имеет многомерное нормальное распределение с параметрами \vec{a} и K . Координаты вектора $\vec{\xi}$ независимы тогда и только тогда, когда они не коррелированы, т.е. матрица ковариаций K диагональная.

Замечание 1.4.19. Если плотность совместного распределения нормальных случайных величин ξ и η ненулевая, то они независимы тогда и только тогда, когда их коэффициент корреляции равен нулю.

Теорема 1.4.20. (Многомерная центральная предельная теорема) Среднее арифметическое независимых одинаково распределенных случайных векторов слабо сходится к многомерному нормальному распределению.

Лемма Фишера

Лм 1.4.21 (Фишера). Пусть \vec{X} состоит из независимых нормальных стандартных величин, а $\vec{Y} = C\vec{X}$, где C — ортогональная матрица. Тогда случайная величина

$$T(\vec{X}) = \sum_{i=1}^n X_i^2 - \sum_{i=1}^k Y_i^2 \quad \forall 1 \leq k \leq n-1$$

не зависит от случайных величин Y_1, \dots, Y_k и имеет распределение H_{n-k} .

□ Т.к. C ортогональная матрица, то $\|\vec{X}\| = \|\vec{Y}\|$, значит

$$\sum_{i=1}^n X_i^2 = \sum_{i=1}^n Y_i^2 \implies T(\vec{X}) = \sum_{i=1}^n Y_i^2 - \sum_{i=1}^k Y_i^2 = \sum_{i=k+1}^n Y_i^2$$

Таким образом $T(\vec{X}) \in H_{n-k}$, т.к. по 1.4.16 случайные величины Y_i независимы и имеют стандартное нормальное распределение. ■

Основная теорема

Замечание 1.4.22. Название неофициальное и актуально только в рамках курса, т.к. мы будем часто ссылаться на эту теорему.

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ из $N(a; \sigma^2)$. Обозначим \bar{x} — выборочное среднее, а S^2 — исправленную выборочную дисперсию. Тогда имеют место следующие утверждения.

Лм 1.4.23.

$$\sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} \in N(0; 1)$$

□

$$X \in N(a; \sigma^2) \implies \sum_{i=1}^n X_i \in N(na; n\sigma^2) \implies \bar{x} \in N\left(a; \frac{\sigma^2}{n}\right) \implies \bar{x} - a \in N\left(0; \frac{\sigma^2}{n}\right) \implies \sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} \in N(0; 1)$$

■

Lm 1.4.24.

$$\sum_{i=1}^n \left(\frac{x_i - a}{\sigma} \right)^2 \in H_n$$

□

$$\forall i \left| \frac{x_i - a}{\sigma} \in N(0; 1) \implies \sum_{i=1}^n \left(\frac{x_i - a}{\sigma} \right)^2 \in H_n \text{ по определению}$$

■

Lm 1.4.25.

$$\sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 = \frac{(n-1)S^2}{\sigma^2} \in H_{n-1}$$

□

$$\sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 = \sum_{i=1}^n \left(\frac{x_i - a}{\sigma} - \frac{\bar{x} - a}{\sigma} \right)^2 = \sum_{i=1}^n (z_i - \bar{z})^2$$

где

$$z_i = \frac{x_i - a}{\sigma} \in N(0; 1)$$

$$\bar{z} = \frac{z_1 + \dots + z_n}{n} = \frac{1}{n} \sum_{i=1}^n \frac{x_i - a}{\sigma} = \frac{\sum x_i - na}{n\sigma} = \frac{n\bar{x} - na}{n\sigma} = \frac{\bar{x} - a}{\sigma}$$

Применим 1.4.21. Сначала рассмотрим

$$T \left(\begin{pmatrix} \vec{z} \\ \bar{z} \end{pmatrix} \right) = \sum_{i=1}^n (z_i - \bar{z})^2 = n\mathbb{D}^* = \sum_{i=1}^n z_i^2 - n(\bar{z})^2 = \sum_{i=1}^n z_i^2 - (\sqrt{n}\bar{z})^2$$

Рассмотрим второе слагаемое.

$$\sqrt{n} \cdot \bar{z} = \frac{z_1 + \dots + z_n}{\sqrt{n}} = \frac{1}{\sqrt{n}} z_1 + \dots + \frac{1}{\sqrt{n}} z_n$$

Строка $\left(\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}} \right)$ имеет единичную длину, поэтому, как известно из курса линейной алгебры, мы можем ее дополнить до ортогональной матрицы C . Тогда $\sqrt{n} \cdot \bar{z} = Y_1$ — первая компонента вектора $\vec{Y} = C\vec{Z}$. Отсюда по 1.4.21 имеем

$$T \left(\begin{pmatrix} \vec{z} \\ \bar{z} \end{pmatrix} \right) = \sum_{i=1}^n (z_i - \bar{z})^2 \in H_{n-1}$$

Причем он не зависит от $Y_1 = \sqrt{n} \cdot \bar{z}$.

■

Lm 1.4.26.

$$\sqrt{n} \cdot \frac{\bar{x} - a}{S} \in T_{n-1}$$

□

$$\sqrt{n} \cdot \frac{\bar{x} - a}{S} = \sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} \cdot \frac{1}{\sqrt{\frac{(n-1)S^2}{\sigma^2} \cdot \frac{1}{n-1}}} = \frac{X_0}{\sqrt{\frac{\chi_{n-1}^2}{n-1}}}$$

где $X_0 \in N(0; 1)$ по 1.4.23, а $\frac{(n-1)S^2}{\sigma^2} \in H_{n-1}$ по 1.4.25. Итого по определению получаем распределение Стьюдента с $n-1$ степенями свободы. Стоит отметить, что в определении распределения Стьюдента числитель и знаменатель должны быть независимы. У нас с этим проблем нет, т.к. \bar{x} и S^2 независимы согласно 1.4.27. ■

Lm 1.4.27.

\bar{x}, S^2 независимые случайные величины

□ В 1.4.25 показали, что

$$\sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 = \frac{(n-1)S^2}{\sigma^2} \text{ не зависит от } \sqrt{n} \cdot \bar{x}$$

Умножение на константу не влияет на независимость, поэтому получаем, что S^2 не зависит от \bar{x} . ■

1.5. Лекция 24.03.07.

Квантили распределений

Предполагаем, что распределение абсолютно непрерывно и $F(X)$ — его функция распределения.

Def 1.5.1 (A). Число t_γ называется квантилем распределения уровня γ , если $F(t_\gamma) = \gamma$. Таким образом $t_\gamma = F^{-1}(\gamma)$.

Замечание 1.5.2. Медиана это квантиль уровня $\frac{1}{2}$.

TODO: image 4:30

В ряде источников под квантилем понимается несколько иное.

Def 1.5.3 (B). Число t_α называется квантилем уровня значимости α , если $P(x > t_\alpha) = \alpha$ или $F(t_\alpha) = 1 - \alpha$. Таким образом t_α соответствует вероятности попадания величины в правостороннюю область.

Замечание 1.5.4. Будем использовать оба определения, но во втором случае добавлять «уровня значимости» и использовать символ α , а не γ .

Квантили основных распределений в Excel

1. НОРМ.СТ.ОБР — обратная функция к функции стандартного нормального распределения.

$$F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{z^2}{2}\right) dz$$

Тогда НОРМ.СТ.ОБР $(x + 0.5)$ это обратная функция к функции Лапласа $\Phi(x)$.

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x \exp\left(-\frac{z^2}{2}\right) dz$$

2. (a) СТЬЮДЕНТ.ОБР (γ, n) — функция, обратная к функции распределения Стьюдента с n степенями свободы.

$$t_n = \frac{X_0}{\sqrt{\frac{1}{n} \chi_n^2}}$$

Заметим, что это распределение симметричное относительно нулю, поэтому можно использовать другой способ.

- (b) СТЬЮДЕНТ.ОБР.2Х — так называемое двустороннее обратное преобразование Стьюдента, которое возвращает t_α такое, что $P(|X| > t_\alpha) = \alpha$.

Таким образом $P(|X| < t_\alpha) = 1 - \alpha$, значит $t_\gamma = \text{СТЮДЕНТ.ОБР.2Х}(1 - \gamma, n)$.

3. (a) ХИ2.ОБР — возвращает квантиль t_γ для распределения H_n с n степенями свободы.
(b) ХИ2.ОБР.ПХ — возвращает квантиль t_α . Заметим, что $\text{ХИ2.ОБР}(\gamma, n) = \text{ХИ2.ОБР.ПХ}(1 - \gamma, n)$.
4. (a) F.ОБР — возвращает t_γ для F -распределения.
(b) F.ОБР.ПХ — возвращает t_α . Заметим, что $\text{F.ОБР}(\gamma, n, m) = \text{F.ОБР.ПХ}(1 - \gamma, n, m)$.

Интервальные оценки

Def 1.5.5. Интервал $(\Theta_\gamma^-; \Theta_\gamma^+)$ называется доверительным интервалом для параметра Θ уровня надежности γ если

$$\Theta \in (\Theta_\gamma^-; \Theta_\gamma^+) \text{ с вероятностью } \gamma \quad \text{или} \quad P(\Theta_\gamma^- < \Theta < \Theta_\gamma^+) = \gamma$$

Замечание 1.5.6. В случае дискретного распределения более точным будет определение

$$P(\Theta_\gamma^- < \Theta < \Theta_\gamma^+) \geq \gamma$$

Замечание 1.5.7. Неверно говорить, что Θ лежит в данном интервале с вероятностью γ , т.к. параметр Θ не является случайной величиной. Правильнее говорить, что доверительный интервал покрывает параметр Θ с вероятностью γ .

Замечание 1.5.8. Вероятность $\alpha = 1 - \gamma$ называется уровнем значимости доверительного интервала.

Замечание 1.5.9. Обычно пытаемся получить доверительный интервал симметричный относительно несмещенной точной оценки Θ^* , но не всегда это возможно.

Замечание 1.5.10. Стандартные уровни надежности 0.9, 0.95, 0.99 и 0.999. Наиболее распространенный 0.95.

Lm 1.5.11. Если случайная величина ξ имеет симметричное распределение относительно точки $x = 0$, то $P(|\xi| < t) = 2F(t) - 1$.

□ Под симметричным распределением подразумевается распределение с симметричным графиком плотности. Тогда

$$P(|\xi| < t) = 2P(0 < \xi < t) = 2(F(t) - F(0)) = 2F(t) - 1$$

т.к. $F(0) = 0.5$ в силу симметрии. ■

Доверительные интервалы для параметров нормального распределения

Пусть $\vec{X} = (X_1, \dots, X_n)$ — выборка объема n из $N(a; \sigma^2)$. Необходимо построить доверительные интервалы. Возможны 4 ситуации.

I. Доверительный интервал для параметра a при известном значении параметра σ^2

По 1.4.23 имеем

$$\sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} \in N(0; 1)$$

Будем искать интервал в виде

$$P\left(-t_\gamma < \sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} < t_\gamma\right) = P\left(\left|\sqrt{n} \cdot \frac{\bar{x} - a}{\sigma}\right| < t_\gamma\right) = 2\Phi(t_\gamma) - 1 = \gamma \implies \Phi(t_\gamma) = \frac{1+\gamma}{2}$$

Таким образом t_γ — обратное значение функции Лапласа или квантиль уровня $\frac{1+\gamma}{2}$ стандартного нормального распределения. Остается решить изначальное неравенство относительно параметра a , получаем

$$\begin{aligned} -t_\gamma &< \sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} < t_\gamma \\ -t_\gamma \frac{\sigma}{\sqrt{n}} &< \bar{x} - a < t_\gamma \frac{\sigma}{\sqrt{n}} \\ \bar{x} - t_\gamma \frac{\sigma}{\sqrt{n}} &< a < \bar{x} + t_\gamma \frac{\sigma}{\sqrt{n}} \end{aligned}$$

Итак, получили доверительный интервал для параметра a надежности γ вида $\left(\bar{x} - t_\gamma \frac{\sigma}{\sqrt{n}}; \bar{x} + t_\gamma \frac{\sigma}{\sqrt{n}}\right)$, где t_γ находится из условия $\Phi(t_\gamma) = \frac{1+\gamma}{2}$.

II. Доверительный интервал для параметра a при неизвестном значении параметра σ^2

По 1.4.26 имеем

$$\sqrt{n} \cdot \frac{\bar{x} - a}{S} \in T_{n-1}$$

Проведем аналогичные рассуждения, что и в прошлом пункте.

$$P\left(-t_\gamma < \sqrt{n} \cdot \frac{\bar{x} - a}{S} < t_\gamma\right) = P\left(\left|\sqrt{n} \cdot \frac{\bar{x} - a}{S}\right| < t_\gamma\right) = 2F_{T_{n-1}}(t_\gamma) - 1 = \gamma \implies F_{T_{n-1}}(t_\gamma) = \frac{1+\gamma}{2}$$

Таким образом t_γ это квантиль распределения Стьюдента уровня $\frac{1+\gamma}{2}$. Решаем неравенство и получаем

$$\bar{x} - t_\gamma \frac{S}{\sqrt{n}} < a < \bar{x} + t_\gamma \frac{S}{\sqrt{n}}$$

Замечание 1.5.12. Заметим, что в обоих случаях получили симметричные интервалы относительно точечной оценки \bar{x} .

Замечание 1.5.13. В Excel квантиль t_γ удобнее находить как $t_\gamma = \text{СТЮДЕНТ.ОБ.2X}(1 - \gamma, n - 1)$.

III. Доверительный интервал для параметра σ^2 при известном значении параметра a

По 1.4.25 имеем

$$\frac{(n-1)S^2}{\sigma^2} \in H_{n-1}$$

Пусть χ_1^2 и χ_2^2 это квантили распределения H_{n-1} уровней $\frac{1-\gamma}{2}$ и $\frac{1+\gamma}{2}$. Тогда

$$P\left(\chi_1^2 < \frac{(n-1)S^2}{\sigma^2} < \chi_2^2\right) = F_{H_{n-1}}(\chi_2^2) - F_{H_{n-1}}(\chi_1^2) = \frac{1+\gamma}{2} - \frac{1-\gamma}{2} = \gamma$$

Решим неравенство относительно σ^2 .

$$\begin{aligned}\chi_1^2 &< \frac{(n-1)S^2}{\sigma^2} < \chi_2^2 \\ \frac{1}{\chi_2^2} &< \frac{\sigma^2}{(n-1)S^2} < \frac{1}{\chi_1^2} \\ \frac{(n-1)S^2}{\chi_2^2} &< \sigma^2 < \frac{(n-1)S^2}{\chi_1^2}\end{aligned}$$

Получили доверительный интервал $\left(\frac{(n-1)S^2}{\chi_2^2}; \frac{(n-1)S^2}{\chi_1^2}\right)$, где χ_1^2 и χ_2^2 это квантили распределения H_{n-1} уровней $\frac{1-\gamma}{2}$ и $\frac{1+\gamma}{2}$.

Замечание 1.5.14. Получили доверительный интервал для среднего квадратического отклонения $\left(\frac{S\sqrt{n-1}}{\chi_2}; \frac{S\sqrt{n-1}}{\chi_1}\right)$.

IV. Доверительный интервал для параметра σ^2 при известном значении параметра a

По 1.4.24 имеем

$$\sum_{i=1}^n \left(\frac{x_i - a}{\sigma}\right)^2 = \frac{n(\sigma^2)^*}{\sigma^2} \in H_n$$

где $(\sigma^2)^* = \frac{1}{n} \sum_{i=1}^n (x_i - a)^2$. Пусть χ_1^2 и χ_2^2 это квантили распределения H_n уровней $\frac{1-\gamma}{2}$ и $\frac{1+\gamma}{2}$. Тогда

$$P\left(\chi_1^2 < \frac{n(\sigma^2)^*}{\sigma^2} < \chi_2^2\right) = F_{H_n}(\chi_2^2) - F_{H_n}(\chi_1^2) = \frac{1+\gamma}{2} - \frac{1-\gamma}{2} = \gamma$$

Аналогично третьему пункту решим неравенство и получим

$$\frac{n(\sigma^2)^*}{\chi_2^2} < \sigma^2 < \frac{n(\sigma^2)^*}{\chi_1^2}$$

Итак, доверительный для σ^2 имеет вид $\left(\frac{n(\sigma^2)^*}{\chi_2^2}; \frac{n(\sigma^2)^*}{\chi_1^2}\right)$, где χ_1^2 и χ_2^2 это квантили распределения H_n уровней $\frac{1-\gamma}{2}$ и $\frac{1+\gamma}{2}$.

Замечание 1.5.15. Получили доверительный интервал для σ вида $\left(\frac{\sqrt{n(\sigma^2)^*}}{\chi_2}; \frac{\sqrt{n(\sigma^2)^*}}{\chi_1}\right)$.

Замечание 1.5.16. Доверительные интервалы полученные в третьем и четвертом пунктах не являются симметричными относительно S^2 , но есть формулы для симметричных доверительных интервалов.

Пример 1.5.17. Пусть $X \in N(a; \sigma^2)$ причем известно, что $\sigma = 3$. В результате обработки выборки объема $n = 36$ получили $\bar{x} = 4.1$. Найти доверительный интервал для параметра a надежности $\gamma = 0.95$.

Решение: Найдем t_γ .

$$2\Phi(t_\gamma) = 0.95 \implies \Phi(t_\gamma) = 0.475 \implies t_\gamma = 1.96$$

Применим формулу из первого пункта.

$$\begin{aligned}\bar{x} - t_\gamma \frac{S}{\sqrt{n}} &< a < \bar{x} + t_\gamma \frac{S}{\sqrt{n}} \\ 4.1 - 1.96 \cdot \frac{3}{\sqrt{36}} &< a < 4.1 + 1.96 \cdot \frac{3}{\sqrt{36}} \\ 3.12 &< a < 5.08\end{aligned}$$

Пример 1.5.18. Пусть $X \in N(a; \sigma^2)$. В результате обработки выборки объема $n = 25$ получили $\bar{x} = 42.32$ и $S = 6.4$. Найти доверительный интервал для параметра a при уровне надежности $\gamma = 0.95$.

Решение: Найдем t_γ . Получим СТЬЮДЕНТ.ОБР.2Х $(1 - 0.95, 25 - 1) = 2.064$. Подставим в формулу из второго пункта.

$$\begin{aligned}\bar{x} - t_\gamma \frac{S}{\sqrt{n}} &< a < \bar{x} + t_\gamma \frac{S}{\sqrt{n}} \\ 42.32 - 2.064 \cdot \frac{6.4}{\sqrt{25}} &< a < 42.32 + 2.064 \cdot \frac{6.4}{\sqrt{25}} \\ 39.678 &< a < 44.962\end{aligned}$$

1.6. Лекция 24.03.14.

Проверка статистических гипотез

Def 1.6.1. Гипотезой H называется предположение о распределении случайной величины.

Def 1.6.2. Гипотеза называется простой, если она однозначно определяет распределение, т.е. $H: \mathcal{F} \in \mathcal{F}_1$, где \mathcal{F}_1 — распределение известного типа с известными параметрами. Все остальные гипотезы называются сложными. Они состоят из конечного или бесконечного числа простых гипотез.

Будем изучать лишь самую простую схему из двух гипотез, где H_0 — это основная (нулевая) гипотеза, а H_1 это альтернативная (конкурирующая) гипотеза, состоящая в том, что основная гипотеза неверна. Стоит отметить, что иногда рассматриваются более сложные схемы, состоящие из набора нескольких гипотез.

Основная гипотеза H_0 принимается или отклоняется при помощи статистики критерия K , где $K(X_1, \dots, X_n) \rightarrow \mathbb{R} = \bar{S} \cup S$. Принимаем гипотезу H_0 , если статистика критерия попала в область \bar{S} , в противном случае принимаем гипотезу H_1 .

Def 1.6.3. Область S называется критической областью критерия.

Def 1.6.4. Точки $t_{кр}$ на границе двух областей называются критическими.

Def 1.6.5. Ошибка первого рода состоит в том, что нулевая гипотеза H_0 отклоняется, хотя она верна.

Def 1.6.6. Ошибка второго рода состоит в том, что отклоняется альтернативная гипотеза H_1 , хотя она верна.

Def 1.6.7. Вероятность α ошибки первого рода называется уровнем значимости критерия.

Def 1.6.8. Вероятность ошибки второго рода обозначается β . Мощностью критерия называется вероятность $1 - \beta$, т.е. вероятность не допущения ошибки второго рода.

Замечание 1.6.9. Заметим, что вероятности α и β не связаны каким-либо соотношением.

Замечание 1.6.10. Ясно, что критерий будет тем лучше, чем меньше вероятности ошибок α и β . При увеличении объема выборки обе эти вероятности уменьшаются, однако при фиксированном объеме выборки попытки уменьшить одну вероятность ведут к увеличению другой.

Способы сравнения критериев

Пусть имеется два критерия K_1 и K_2 с соответствующими вероятностями ошибок α_1, β_1 и α_2, β_2 . Объем выборки n фиксированный.

I. Минимаксный подход

Критерий K_1 не хуже, чем K_2 , если $\max(\alpha_1, \beta_1) < \max(\alpha_2, \beta_2)$. Из всех критериев выбираем тот, который не хуже остальных.

II. Байесовский подход

Пусть известны потери h_1 и h_2 от ошибок первого и второго рода. Тогда средние ожидаемые потери составят $u = \alpha h_1 + \beta h_2$. Из всех критериев выбираем тот, у которого средние ожидаемые потери наименьшие.

III. Выбор наиболее мощного критерия

Обозначим $K_\varepsilon = \{K_i \mid \alpha \leq \varepsilon\}$, т.е. класс критериев данной выборки.

Def 1.6.11. Критерий $K \in K_\varepsilon$ называется наиболее мощным критерием уровня ε , если $\beta \leq \beta_i \mid \forall K_i \in K_\varepsilon$.

Построение критериев согласия

Def 1.6.12. Говорят, что критерий K является критерием асимптотического уровня ε , если $\alpha \rightarrow \varepsilon$ при $n \rightarrow \infty$.

Def 1.6.13. Критерий K для проверки гипотезы H_0 против альтернативной гипотезы H_1 называется состоятельным, если $\beta \rightarrow 0$ при $n \rightarrow \infty$.

Def 1.6.14. Критерием согласия уровня ε называется состоятельный критерий асимптотического уровня ε .

В качестве критерия согласия берется статистика $K(X_1, \dots, X_n)$ со свойствами

1. Если H_0 верна, то $K(X_1, \dots, X_n) \rightrightarrows z$, где z — случайная величина с известным распределением.
2. Если H_0 неверна, то $K(X_1, \dots, X_n) \xrightarrow{P} \infty$ при $n \rightarrow \infty$.

Для заданного уровня значимости ε находим квантиль $t_{кр}$ такой, что $P(|z| \geq t_{кр}) = \varepsilon$. В результате получаем критерий с уровнем значимости α . Итого

$$\begin{cases} H_0, & \text{если } |K| < t_{кр} \\ H_1, & \text{если } |K| \geq t_{кр} \end{cases}$$

Замечание 1.6.15. Таким образом в качестве статистики берется функция отклонения эмпирического распределения от теоретического. Если нулевая гипотеза верна, то эта функция сходится к собственному распределению, если нет, то неограниченно возрастает.

Лм 1.6.16. Построенный критерий обладает свойствами

1. Это критерий асимптотического уровня ε .
2. Этот критерий состоятельный.

□ Пусть гипотеза H_0 верна. Тогда по первому свойству выбранной статистики имеем $\forall x \left| F_K(x) \rightarrow F_z(x) \right|$ при $n \rightarrow \infty$. Значит

$$\alpha = P(|K| \geq t_{кр} | H_0) = 1 - (F_K(t_{кр}) - F_K(-t_{кр})) \xrightarrow{n \rightarrow \infty} 1 - (F_z(t_{кр}) - F_z(-t_{кр})) = P(|z| \geq t_{кр}) = \varepsilon$$

Далее покажем состоятельность критерия. По второму свойству выбранной статистики получаем, что если H_1 верна, то $|K| \xrightarrow{P} \infty$, т.е. $\forall c \in \mathbb{R} \left| P(|K| \geq c) \rightarrow 1 \right|$, значит $P(|K| < c) \rightarrow 0$. ■

Гипотеза о среднем нормальной совокупности с известной дисперсией

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ объема n случайной величины $X \in N(a; \sigma^2)$, где σ^2 это известное значение. Проверяется гипотеза H_0 , которая состоит в том, что $a = a_0$. Альтернативная гипотеза H_1 состоит в том, что $a \neq a_0$. В качестве статистики критерия возьмем функцию $K = \sqrt{n} \cdot \frac{\bar{x} - a_0}{\sigma}$. Проверим, что она имеет требуемые свойства.

1. Если H_0 верна, то

$$K = \sqrt{n} \cdot \frac{\bar{x} - a_0}{\sigma} = \sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} \in N(0; 1)$$

по 1.4.23.

2. Если H_0 неверна, т.е. $a \neq a_0$, то

$$|K| = \left| \sqrt{n} \cdot \frac{\bar{x} - a_0}{\sigma} \right| = \left| \sqrt{n} \cdot \frac{\bar{x} - a}{\sigma} + \sqrt{n} \cdot \frac{a - a_0}{\sigma} \right|$$

Первое слагаемое по 1.4.23 имеет стандартное нормальное распределение, поэтому оно будет ограничено по вероятности. Второе слагаемое будет стремиться к бесконечности, т.к. $\sqrt{n} \rightarrow \infty$ при $n \rightarrow \infty$, а оставшая часть этого слагаемого это какая ненулевая константа в силу того, что $a \neq a_0$. Итого $|K| \xrightarrow{P} \infty$ при $n \rightarrow \infty$.

Итак, получили следующий критерий: для уровня значимости α выберем $t_{кр}$ такую, что

$$\alpha = P(|K| \geq t_{кр}) = 1 - 2\Phi(t_{кр}) \implies \Phi(t_{кр}) = \frac{1 - \alpha}{2}$$

Значит $t_{кр}$ это обратное значение функции Лапласа в точке $\frac{1-\alpha}{2}$ или квантиль уровня $1 - \frac{\alpha}{2}$ стандартного нормального распределения. Получаем критерий согласия

$$\begin{cases} H_0, & |K| < t_{кр} \\ H_1, & |K| \geq t_{кр} \end{cases}$$

Гипотеза о среднем нормальной совокупности при неизвестной дисперсии

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ объема n случайной величины $X \in N(a; \sigma^2)$. Проверяется гипотеза H_0 , которая состоит в том, что $a = a_0$. Альтернативная гипотеза H_1 состоит в том, что $a \neq a_0$.

В качестве статистики критерия возьмем функцию $K = \sqrt{n} \cdot \frac{\bar{x} - a_0}{S}$. Проверим, что она имеет требуемые свойства.

1. Если H_0 верна, то

$$K = \sqrt{n} \cdot \frac{\bar{x} - a_0}{S} = \sqrt{n} \cdot \frac{\bar{x} - a}{S} \in T_{n-1}$$

по 1.4.26.

2. Если H_0 неверна, т.е. $a \neq a_0$, то

$$|K| = \left| \sqrt{n} \cdot \frac{\bar{x} - a_0}{S} \right| = \left| \sqrt{n} \cdot \frac{\bar{x} - a}{S} + \sqrt{n} \cdot \frac{a - a_0}{S} \right|$$

Аналогично предыдущей гипотезе получаем, что $|K| \xrightarrow{P} \infty$ при $n \rightarrow \infty$.

Итак, получили следующий критерий: для уровня значимости α выберем $t_{кр}$ такую, что $P(|t_{n-1}| \geq t_{кр}) = \alpha$. Таким образом $t_{кр}$ это квантиль уровня $1 - \frac{\alpha}{2}$ распределения T_{n-1} или $t_{кр}$ это квантиль уровня значимости α двустороннего распределения $|T_{n-1}|$. Получаем критерий согласия

$$\begin{cases} H_0, & |K| < t_{кр} \\ H_1, & |K| \geq t_{кр} \end{cases}$$

Доверительные интервалы как критерии гипотез по параметрам распределения

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ объема n из $X \in \mathcal{F}_\Theta$, где \mathcal{F}_Θ это распределение известного типа, с неизвестным параметром Θ . Проверяется гипотеза H_0 о том, что $\Theta = \Theta_0$, против альтернативной гипотезы $\Theta \neq \Theta_0$. Допустим, что для параметра Θ построен доверительный интервал $(\Theta_\gamma^-; \Theta_\gamma^+)$ надежности γ .

Лм 1.6.17. В описанных условиях получаем критерий

$$\begin{cases} H_0, & \Theta_0 \in (\Theta_\gamma^-; \Theta_\gamma^+) \\ H_1, & \Theta_0 \notin (\Theta_\gamma^-; \Theta_\gamma^+) \end{cases}$$

уровня $\alpha = 1 - \gamma$.

□

$$\alpha = P(\Theta_0 \notin (\Theta_\gamma^-; \Theta_\gamma^+) \mid X \in \mathcal{F}_{\Theta_0}) = 1 - P(\Theta_0 \in (\Theta_\gamma^-; \Theta_\gamma^+) \mid X \in \mathcal{F}_{\Theta_0}) = 1 - \gamma$$

■

Замечание 1.6.18. Доказать, что данный критерий будет состоятельным, в общем случае нельзя.

Пример 1.6.19. По выборке объема $n = 36$ из нормальной совокупности с известным средним квадратическим отклонением $\sigma = 1.44$ найдено выборочное среднее $\bar{x} = 21.6$. Проверить гипотезу H_0 о том, что $a = 21$, против альтернативной гипотезы H_1 о том, что $a \neq 21$. Уровень значимости принять $\alpha = 0.05$.

Решение: Составляем статистику

$$K = \sqrt{n} \cdot \frac{\bar{x} - a_0}{\sigma} = \sqrt{36} \cdot \frac{21.6 - 21}{1.44} = 2.5$$

Теперь ищем критическую точку.

$$\Phi(t_{кр}) = \frac{1 - \alpha}{2} = 0.475 \implies t_{кр} = 1.96$$

Т.к. $|K| > t_{кр}$, то основная гипотеза H_0 отклоняется, и принимается альтернативная гипотеза H_1 .

Некоторые дополнения с практики

Теорема 1.6.20. Если $f_\xi(Me) \neq 0$, то $Me^* \xrightarrow{P} Me$, причем со скоростью $\frac{1}{\sqrt{n}}$. Таким образом выборочная медиана является состоятельной асимптотически нормальной оценкой теоретической медианы.

В случае симметричных распределений выборочную медиану удобно использовать, если распределение имеет «жирные хвосты». В случае «жирных» хвостов возможны так называемые «выбросы», которые заметно повлияют на выборочное среднее, и поэтому его не желательно использовать в качестве оценки центра симметрии. В случае «жирных хвостов» или «выбросов» можно использовать медиану или иные методики.

Недостатки медианы:

1. Хорошо работает только в случае симметричных распределений.
2. В случае нормального распределения медиана сходится приблизительно на 20% медленнее, чем среднее выборочное, и до 40% медленнее в некоторых других случаях.

Подумаем над устранением этих недостатков. Пусть распределение симметрично. Рассмотрим два метода.

I. Метод усеченного среднего

Метод заключается в том, что мы отбрасываем из выборки сверху и снизу по 5–10 значений или по 5% значений.

Пусть дана выборка $\vec{X} = (X_1, \dots, X_n)$ объема n . Допустим, что мы отбросили по k значений сверху и снизу и получили выборку объема $n - 2k$ вида $(X_{(k+1)}, \dots, X_{(n-k)})$. Вычислим среднее выборочное этой выборки $\frac{\sum_{i=k+1}^{n-k} X_i}{n - 2k}$. Это компромисс между средним выборочным и медианой, т.к.

1. При $k = 0$ получаем \bar{x} .
2. При $n - 2k = 1$ получаем Me .

II. Метод средних Уолша

Пусть дана выборка $\vec{X} = (X_1, \dots, X_n)$ объема n . Для всех пар считаем среднее арифметическое $y_k = \frac{1}{2}(x_i + x_j)$, где $1 \leq k \leq \frac{1}{2}n(n+1)$. Далее находим медиану полученной новой выборки. Преимущества данного метода:

1. Эффективность по сравнению со средним выборочным упадет не более, чем на 12%.
2. Сглаживает выбросы.

1.7. Лекция 24.03.21.

Критерии для проверки гипотез о распределении

I. Критерий χ^2 (для проверки параметрической гипотезы)

Пусть дана выборка $\vec{X} = (X_1, \dots, X_n)$ неизвестного распределения \mathcal{F} . Проверяется основная (сложная) гипотеза H_0 о том, что $\mathcal{F} \in \mathcal{F}_\Theta$, где \mathcal{F}_Θ это распределение известного типа, определяемое набором из m неизвестных параметров $\Theta = (\Theta_1, \dots, \Theta_m)$. В качестве альтернативной гипотезы H_1 о том, что $\mathcal{F} \notin \mathcal{F}_\Theta$.

Пусть $\hat{\Theta} = (\hat{\Theta}_1, \dots, \hat{\Theta}_m)$ — оценка неизвестных параметров $\Theta = (\Theta_1, \dots, \Theta_m)$ методом максимального правдоподобия.

Разобьем выборку на k интервалов $A_i = [a_i; a_{i+1})$, обозначим n_i соответствующие частоты этих интервалов. Пусть $p_i = F_{\hat{\Theta}}(a_{i+1}) - F_{\hat{\Theta}}(a_i)$ это теоретические вероятности попадания случайной величины с распределением $\mathcal{F}_{\hat{\Theta}}$ в данные интервалы. Тогда $n'_i = np_i$ это теоретические частоты в эти интервалы.

В качестве статистики критерия берется функция

$$K = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i} = \sum_{i=1}^k \frac{n_i^2}{n'_i} - n$$

Теорема 1.7.1. (Фишера) Если гипотеза $H_0: \mathcal{F} \in \mathcal{F}_\Theta$ верна, то

$$K = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i} \Rightarrow \chi^2_{k-m-1}$$

где k это число интервалов, m — число параметров распределения.

Получили следующий критерий: для заданного уровня значимости α находим $t_{кр}$ такое, что $P(\chi^2_{k-m-1} \geq t_{кр}) = \alpha$. Тогда критерий согласия имеет вид

$$\begin{cases} H_0, & K < t_{кр} \\ H_1, & K \geq t_{кр} \end{cases}$$

Замечание 1.7.2. Для вычисления квантиля можно использовать формулу Excel **ХИ2.ОБР.ПХ** ($\alpha, k - m - 1$).

Замечание 1.7.3. Частота каждого интервала должна быть не менее пяти. Если это не так, то сливаем соседние интервалы. Число интервалов желательно брать достаточно большим (чтобы $k - m - 1$ было больше), но при этом стоит помнить об ограничениях на частоту каждого интервала.

Замечание 1.7.4. При этом критерии выборку лучше разбивать на «равнонаполненные» интервалы, т.е. на интервалы содержащие примерно одинаковое число элементов выборки.

Замечание 1.7.5. Желательно, чтобы объем выборки был не менее 50. В противном случае метод работает плохо.

A_i	[5.2; 7.4)	[7.4; 9.6)	[9.6; 11.8)	[11.8; 14.0)	[14.0; 16.2)	[16.2; 18.4)	[18.4; 20.6)	[20.6; 22.8]	\sum
n_i	12	17	14	13	18	14	13	19	120

Таблица 1.7.6: Интервальный ряд для примера 1.7.7

Пример 1.7.7. Имеется выборка в виде вариационного ряда $\vec{X} = (5.2, \dots, 22.8)$ объема $n = 120$. При разбиении ее на $k = 8$ интервалов одинаковой длины получили интервальный ряд 1.7.6. Проверить гипотезу о равномерности этого распределения при уровне значимости $\alpha = 0.05$.

Решение: Итак, проверяем гипотезу $H_0: \mathcal{F} \in U(a; b)$ против альтернативной $H_1: \mathcal{F} \notin U(a; b)$. По методу наибольшего правдоподобия получаем (смещенные) оценки $a^* = 5.2$ и $b^* = 22.8$. Теоретические вероятности случайной величины с распределением $U(5.2; 22.8)$ равны $p_i = \frac{1}{k} = \frac{1}{8}$. Тогда теоретические частоты $n'_i = np_i = 15$. Далее вычисляем статистику

$$\chi_{\text{набл}}^2 = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i} = 3.2$$

При уровне значимости $\alpha = 0.05$ и числе степеней свободы $k - m - 1 = 5$ находим $t_{\text{кр}} = \text{ХИ2.ОБР.ПХ}(0.05, 5) = 11.07$. Таким образом, т.к. $\chi_{\text{набл}} < t_{\text{кр}}$, то принимаем основную гипотезу о равномерном распределении.

II. Критерий χ^2

Постановка задачи отличается от предыдущего пункта тем, что проверяется простая основная гипотеза H_0 о том, что $\mathcal{F} = \mathcal{F}_1$, где \mathcal{F}_1 это распределение известного типа с известными параметрами. В качестве статистики берем ту же функцию

$$K = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i}$$

Теорема 1.7.8. (Пирсона) Если $H_0: \mathcal{F} = \mathcal{F}_1$ верна, то

$$K = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i} \Rightarrow \chi_{k-1}^2$$

В остальном критерий аналогичен предыдущему.

III. Критерий Колмогорова

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ объема n распределения \mathcal{F} . Проверяется нулевая гипотеза $H_0: \mathcal{F} = \mathcal{F}_1$ против альтернативной $H_1: \mathcal{F} \neq \mathcal{F}_1$, где \mathcal{F}_1 это абсолютно непрерывное распределение известного типа с известными параметрами. Обозначим $F(x)$ — функцию распределения \mathcal{F}_1 , а $F^*(x)$ — выборочную функцию распределения.

Теорема 1.7.9. (Колмогорова) Если гипотеза $H_0: \mathcal{F} = \mathcal{F}_1$ верна, то

$$K = \sqrt{n} \cdot \sup_x |F(x) - F^*(x)| \Rightarrow \mathcal{K}$$

где \mathcal{K} это распределение Колмогорова с функцией распределения

$$F_{\mathcal{K}}(x) = \sum_{j=-\infty}^{\infty} (-1)^j \exp(-2j^2 x^2)$$

В результате получаем критерий Колмогорова: пусть $t_{\text{кр}}$ это квантиль распределения Колмогорова уровня значимости α , тогда

$$\begin{cases} H_0, & K < t_{\text{кр}} \\ H_1, & K \geq t_{\text{кр}} \end{cases}$$

Критерии для проверки однородности

IV. Критерий Колмогорова—Смирнова

Пусть имеются две независимые выборки $\vec{X} = (X_1, \dots, X_n)$ и $\vec{Y} = (Y_1, \dots, Y_m)$ неизвестных непрерывных распределений \mathcal{F} и \mathcal{J} . Проверяется основная гипотеза H_0 о том, что $\mathcal{F} = \mathcal{J}$ против альтернативной H_1 о том, что $\mathcal{F} \neq \mathcal{J}$. В качестве статистики критерия берется функция

$$K = \sqrt{\frac{nm}{n+m}} \cdot \sup_x |\mathcal{F}^*(x) - \mathcal{J}^*(x)|$$

где $\mathcal{F}^*(x)$ и $\mathcal{J}^*(x)$ это соответствующие выборочные функции распределения.

Теорема 1.7.10. (Колмогорова—Смирнова) Если гипотеза $H_0: \mathcal{F} = \mathcal{J}$ верна, то $K \Rightarrow \mathcal{K}$.

В результате получаем критерий Колмогорова—Смирнова: пусть $t_{\text{кр}}$ это квантиль распределения Колмогорова уровня значимости α , тогда

$$\begin{cases} H_0, & K < t_{\text{кр}} \\ H_1, & K \geq t_{\text{кр}} \end{cases}$$

Замечание 1.7.11. Этот критерий не часто используют. Основные критерии это критерии Фишера и Стьюдента (они описаны ниже). Сначала применяется критерий Фишера, и если он не отвергает основную гипотезу, применяется критерий Стьюдента.

V. Критерий Фишера

Пусть имеются две независимые выборки $\vec{X} = (X_1, \dots, X_n)$ и $\vec{Y} = (Y_1, \dots, Y_m)$ нормальных распределений $X \in N(a_1; \sigma_1^2)$ и $Y \in N(a_2; \sigma_2^2)$. Проверяется основная гипотеза H_0 о том, что $\sigma_1 = \sigma_2$ против альтернативной H_1 о том, что $\sigma_1 \neq \sigma_2$. В качестве статистики критерия берется функция

$$K = \frac{S_x^2}{S_y^2}$$

где S_x^2 и S_y^2 — соответствующие исправленные выборочные дисперсии, причем $S_x^2 \geq S_y^2$.

Теорема 1.7.12. Если гипотеза $H_0: \sigma_1 = \sigma_2$ верна, то $K \in F_{n-1, m-1}$.

□ Если H_0 верна, то $\sigma_1 = \sigma_2$, поэтому

$$K = \frac{S_x^2}{S_y^2} = \frac{S_x^2}{\sigma_1^2} \cdot \frac{\sigma_2^2}{S_y^2}$$

По 1.4.25 получаем, что

$$\frac{(n-1)S^2}{\sigma^2} \in H_{n-1} \implies K = \frac{\chi_{n-1}^2}{n-1} \cdot \frac{m-1}{\chi_{m-1}^2} \in F_{n-1, m-1} \text{ по определению}$$

■

Получили критерий: пусть $t_{кр}$ это квантиль распределения $F_{n-1, m-1}$ уровня значимости α , тогда

$$\begin{cases} H_0, & K < t_{кр} \\ H_1, & K \geq t_{кр} \end{cases}$$

Замечание 1.7.13. Если гипотеза неверна, то статистика K стремится не к бесконечности, а к $\frac{\sigma_1^2}{\sigma_2^2}$.

VI. Критерий Стьюдента

Пусть имеются две независимые выборки $\vec{X} = (X_1, \dots, X_n)$ и $\vec{Y} = (Y_1, \dots, Y_m)$ нормальных распределений $X \in N(a_1; \sigma^2)$ и $Y \in N(a_2; \sigma^2)$ с одинаковой дисперсией σ^2 . Проверяется основная гипотеза H_0 о том, что $a_1 = a_2$ против альтернативной H_1 о том, что $a_1 \neq a_2$. В качестве статистики критерия берется функция

$$K = \sqrt{\frac{nm}{n+m}} \cdot \frac{\bar{x} - \bar{y}}{\sqrt{\frac{(n-1)S_x^2 + (m-1)S_y^2}{n+m-2}}}$$

Теорема 1.7.14.

$$\sqrt{\frac{nm}{n+m}} \cdot \frac{(\bar{x} - a_1) - (\bar{y} - a_2)}{\sqrt{\frac{(n-1)S_x^2 + (m-1)S_y^2}{n+m-2}}} \in T_{n+m-2}$$

Из этой теоремы следует, что если гипотеза H_0 верна, то $K \in T_{n+m-2}$. А если она не верна, то $|K| \rightarrow \infty$. Получаем критерий: пусть $t_{кр}$ это квантиль распределения $|T_{n+m-2}|$ уровня значимости α , тогда

$$\begin{cases} H_0, & |K| < t_{кр} \\ H_1, & |K| \geq t_{кр} \end{cases}$$

Критерий вероятности появления события

Пусть $P(A) = p$ — неизвестная теоретическая вероятность события A . При достаточно большом числе n независимых испытаний события A появилось m раз. Проверяется основная гипотеза H_0 о том, что $P(A) = p_0$ против альтернативной H_1 о том, что $P(A) \neq p_0$. Обозначим $p^* = \frac{m}{n}$ точечную оценку $P(A)$. В качестве статистики критерия берем величину

$$K = \frac{m - np_0}{\sqrt{np_0q_0}} \quad (q_0 = 1 - p_0)$$

Если H_0 верна, то $K = \frac{m - np}{\sqrt{npq}} \Rightarrow N(0; 1)$ согласно теореме Муавра—Лапласа. Получаем критерий: пусть $t_{кр}$ это точка такая, что $\Phi(t_{кр}) = \frac{1-\alpha}{2}$, где α это уровень значимости. Тогда

$$\begin{cases} H_0, & |K| < t_{кр} \\ H_1, & |K| \geq t_{кр} \end{cases}$$

Пример 1.7.15. При посеве гороха вероятность рецессивного признака $p = 0.25$, а доминантного — $q = 0.75$. Из $n = 4000$ посеянных семян 970 оказались с рецессивным признаком, а 3030 — с доминантным. Проверяется основная гипотеза H_0 о том, что $p = 0.25$, против альтернативной H_1 о том, что $p \neq 0.25$.

Решение: Вычисляем статистику

$$K = \frac{m - np_0}{\sqrt{np_0q_0}} = -1.095$$

Далее найдем $t_{кр}$.

$$\Phi(t_{кр}) = \frac{1 - 0.05}{2} = 0.475 \implies t_{кр} = 1.96$$

Т.к. $|K| < t_{кр}$, то гипотеза H_0 принимается.

1.8. Лекция 24.03.28.

Статистическая зависимость

Def 1.8.1. Зависимость называется статистической, если изменение одной случайной величины вызывает изменение распределения другой. Если при этом изменяется среднее значение другой случайной величины, то такая статистическая зависимость называется корреляционной. Если при увеличении одной случайной величины, среднее другой также увеличивается, то имеет место прямая корреляция, в противном случае — обратная корреляция.

Корреляционное облако

Пусть в ходе n экспериментов появились значения случайных величин X и Y . Нанеся эти точки на координатную плоскость (рис. 1.8.2) получим корреляционное облако. По его виду можно сделать некоторые предположения о наличии и типе корреляции.

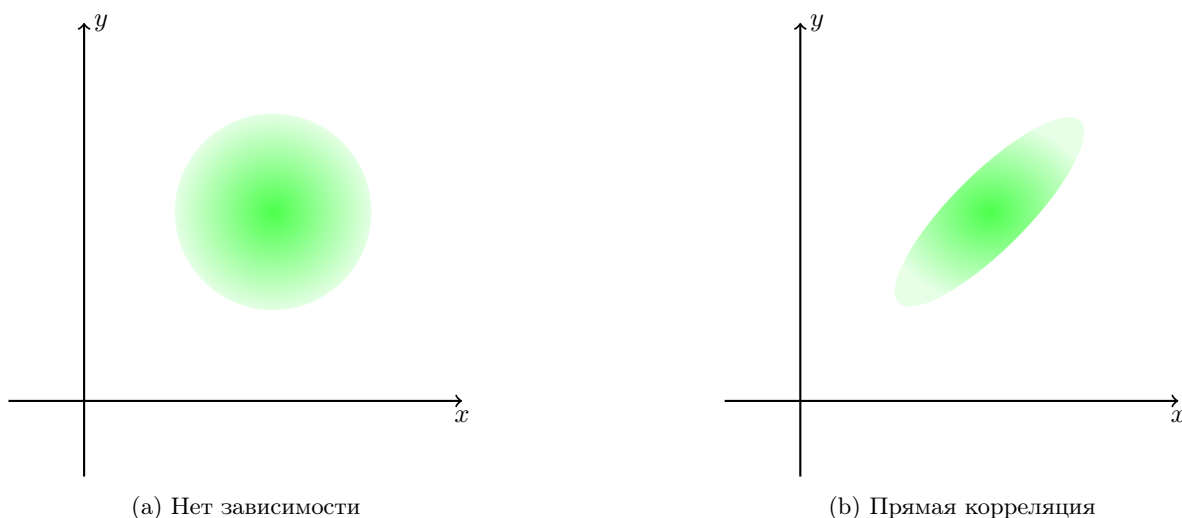


Рис. 1.8.2: Корреляционное облако

Корреляционная таблица

Пусть даны экспериментальные данные $(x_1, y_1), \dots, (x_n, y_n)$. Их удобно представить в виде корреляционной таблицы, где по горизонтали будут расположены x_i , по вертикали — y_i , а в клетках будут находиться числа $n_{i,j}$, показывающие сколько раз встретилась точка (x_i, y_j) .

X \ Y	10	20	30	40	n_x	\bar{y}_x
2	7	3	0	0	10	13
4	3	10	10	2	25	24.4
6	0	2	10	3	15	30.67
n_y	10	15	20	5	$\sum = 50$	

Таблица 1.8.3: Корреляционная таблица

Пример 1.8.4. Пусть $n = 50$. Данные представлены в виде корреляционной таблицы 1.8.3. Столбец \bar{y}_x не входит в корреляционную таблицу, но обычно записывается рядом. Это условное среднее, т.е. выборочное среднее одной случайной величины при условии, что другая величина зафиксирована. Его можно вычислить по формуле

$$\bar{y}_x = \frac{1}{n_x} \sum_i n_{x,y} y_i \quad \bar{x}_y = \frac{1}{n_y} \sum_i n_{x,y} x_i$$

Т.к. с ростом x растут условные средние \bar{y}_x , то имеет место прямая корреляция.

Замечание 1.8.5. Если основные частоты сгруппированы относительно главной диагонали, то имеет место прямая корреляция, а если относительно побочной, то обратная. Если частоты распределены примерно равномерно, то корреляции практически нет.

Замечание 1.8.6. При большом объеме данных непрерывных случайных величин X и Y данные удобно собрать в интервальную корреляционную таблицу, где по вертикали отмечаем интервалы $[a_i; a_{i+1})$ случайной величины X , а по вертикали — интервалы $[b_j; b_{j+1})$ случайной величины Y . В клетках указываем количество точек $v_{i,j}$, попавших в соответствующую прямоугольную область $[a_i; a_{i+1}) \times [b_j; b_{j+1})$.

Критерий χ^2 для проверки независимости

Пусть выборка $(x_1, y_1), \dots, (x_n, y_n)$ представлена в виде интервальной корреляционной таблицы. Случайная величина X при этом разбита на k интервалов, а случайная величина Y на m интервалов. Обозначим $v_{i\cdot}$ число значений X в интервале $[a_i; a_{i+1})$, где $1 \leq i \leq k$, $v_{\cdot j}$ число значений Y в интервале $[b_j; b_{j+1})$, где $1 \leq j \leq m$, а $v_{i,j}$ это число точек в соответствующем прямоугольнике $[a_i; a_{i+1}) \times [b_j; b_{j+1})$. Данные собрали в корреляционную таблицу 1.8.7.

$X \backslash Y$	$[b_0; b_1)$	$[b_1; b_2)$	\dots	$[b_{m-1}; b_m]$	$v_{i\cdot} = \sum_{j=1}^m v_{i,j}$
$[a_0; a_1)$	$v_{1,1}$	$v_{1,2}$	\dots	$v_{1,m}$	$v_{1\cdot}$
\dots	\dots	\dots	\dots	\dots	\dots
$[a_{k-1}; a_k]$	$v_{k,1}$	$v_{k,2}$	\dots	$v_{k,m}$	$v_{k\cdot}$
$v_{\cdot j} = \sum_{i=1}^k v_{i,j}$	$v_{\cdot 1}$	$v_{\cdot 2}$	\dots	$v_{\cdot m}$	n

Таблица 1.8.7: Корреляционная таблица

Проверяется основная гипотеза H_0 о том, что X и Y — независимы против альтернативной H_1 , а том, что они зависимы. Если гипотеза H_0 верна, то теоретическая вероятность попадания случайной величины $\langle x, y \rangle$ в любой прямоугольник равна произведению теоретических вероятностей попадания этих случайных величин в соответствующие интервалы. Таким образом

$$p_{i,j} = P(X \in [a_{i-1}; a_i) \text{ и } Y \in [b_{j-1}; b_j)) = P(X \in [a_{i-1}; a_i)) P(Y \in [b_{j-1}; b_j)) = p_i q_j$$

Тогда по Закону Больших Чисел

$$\frac{v_{i\cdot}}{n} \xrightarrow{P} p_i \quad \frac{v_{\cdot j}}{n} \xrightarrow{P} q_j \quad \frac{v_{i,j}}{n} \xrightarrow{P} p_{i,j}$$

поэтому основанием для отклонения нулевой гипотезы должна служить заметная разница между величинами $\frac{v_{i,j}}{n}$ и $\frac{v_{i\cdot}}{n} \cdot \frac{v_{\cdot j}}{n}$. В качестве статистики критерия берется функция

$$K = n \sum_{i,j} \frac{(v_{i,j} - \frac{1}{n} v_{i\cdot} v_{\cdot j})^2}{v_{i\cdot} v_{\cdot j}}$$

Теорема 1.8.8. Если основная гипотеза H_0 верна, то $K \Rightarrow \chi^2_{(k-1)(m-1)}$.

Получили критерий: пусть $t_{кр}$ это квантиль распределения $H_{(n-1)(m-1)}$ уровня значимости α , тогда

$$\begin{cases} H_0, & K < t_{кр} \\ H_1, & K \geq t_{кр} \end{cases}$$

Замечание 1.8.9. Частота каждой клетки должна быть не менее пяти.

Однофакторный дисперсионный анализ

Предположим, что на случайную величину X (результат) может влиять фактор Z , причем Z не обязательно случайная величина. Требуется определить, оказывает ли фактор Z на среднее значение X . Пусть при различных k уровнях фактора Z получены k независимых выборок случайной величины X . Обозначим их $X^{(1)} = (X_1^{(1)}, \dots, X_{n_1}^{(1)})$, \dots , $X^{(k)} = (X_1^{(k)}, \dots, X_{n_k}^{(k)})$. В общем случае их распределение отличается, поэтому с формальной точки зрения это выборки различных случайных величин.

Общая межгрупповая и внутригрупповая дисперсия

Для каждой выборки вычислим ее выборочное среднее и дисперсию.

$$\bar{x}^{(j)} = \frac{1}{n_j} \sum_{i=1}^{n_j} x_i^{(j)} \quad \mathbb{D}^{(j)} = \frac{1}{n_j} \sum_{i=1}^{n_j} \left(x_i^{(j)} - \bar{x}^{(j)} \right)^2$$

Объединив все выборки в одну общую получаем выборку объема $n = n_1 + \dots + n_k$. Для нее вычислим общее выборочное среднее и общую выборочную дисперсию

$$\bar{x} = \frac{1}{n} \sum_{i,j} x_i^{(j)} = \frac{1}{n} \sum_{j=1}^k \bar{x}^{(j)} n_j \quad \mathbb{D}_O = \frac{1}{n} \sum_{i,j} \left(x_i^{(j)} - \bar{x} \right)^2$$

Def 1.8.10. Внутригрупповой (остаточной) дисперсией называется среднее групповых дисперсий.

$$\mathbb{D}_B = \frac{1}{n} \sum_{j=1}^k n_j \mathbb{D}^{(j)}$$

Def 1.8.11. Межгрупповой (факторной) дисперсией (или дисперсией выборочных средних) называется величина

$$\mathbb{D}_M = \frac{1}{n} \sum_{j=1}^k \left(\bar{x}^{(j)} - \bar{x} \right)^2 n_j$$

Теорема 1.8.12. (О разложении дисперсии)

$$\mathbb{D}_O = \mathbb{D}_M + \mathbb{D}_B$$

□ Это можно доказать чисто алгебраически исходя из определений. ■

Замечание 1.8.13. Ясно, что чем большее влияние оказывает Z на результат X , тем больше получаются отдельные выборочные средние, а значит растет доля межгрупповой дисперсии в данной сумме. Таким образом величина корреляции характеризуется отношением $\frac{\mathbb{D}_M}{\mathbb{D}_O}$. Если зависимость функциональная, то эта дробь будет равна единице.

Проверка гипотезы о влиянии фактора

Предполагаем, что случайная величина X имеет нормальное распределение и фактор Z может влиять только на среднее значение, но не на дисперсию и тип распределения. Из-за этого можно считать, что данные независимые k -выборки при различных уровнях фактора Z также имеют нормальное распределение с одинаковой дисперсией, т.е. $X^{(j)} \in N(a_j; \sigma^2)$.

Проверяется основная гипотеза H_0 о том, что $a_1 = \dots = a_j$, т.е. Z не влияет на среднее X , против альтернативной H_1 о том, что фактор Z влияет на среднее X . По 1.4.25 для каждой из k выборок имеем

$$\sum_{i=1}^n \left(\frac{x - \bar{x}}{\sigma} \right)^2 = \frac{n \mathbb{D}^*}{\sigma^2} \in H_{n_j-1} \quad 1 \leq j \leq k$$

Т.к. распределение «хи-квадрат» устойчиво относительно суммирования, то

$$\sum_{j=1}^k \frac{n_j \mathbb{D}^{(j)}}{\sigma^2} = \frac{n \mathbb{D}_B}{\sigma^2} \in H_{n-k} \text{ т.к. } \underbrace{(n_1 - 1) + \dots + (n_k - 1)}_{k \text{ раз}} = n - k$$

Пусть основная гипотеза H_0 верна, тогда все данные можно считать одной выборкой объема n и опять по 1.4.25 имеем

$$\frac{n \mathbb{D}_O}{\sigma^2} \in H_{n-1}$$

Согласно 1.8.12 получаем

$$\begin{aligned} \mathbb{D}_O &= \mathbb{D}_M + \mathbb{D}_B \\ \frac{n \mathbb{D}_O}{\sigma^2} &= \frac{n \mathbb{D}_M}{\sigma^2} + \frac{n \mathbb{D}_B}{\sigma^2} \end{aligned}$$

Т.к. левая часть имеет распределение H_{n-1} , а второе слагаемое в правой части имеет распределение H_{n-k} , то первое слагаемое в правой части имеет распределение H_{k-1} . Итак, при первой основной гипотезе H_0 получили, что

$$\frac{n \mathbb{D}_M}{\sigma^2} \in H_{k-1} \quad \frac{n \mathbb{D}_B}{\sigma^2} \in H_{n-k}$$

Значит

$$\frac{n\mathbb{D}_M}{\sigma^2(k-1)} \cdot \frac{\sigma^2(n-k)}{n\mathbb{D}_B} = \frac{n-k}{k-1} \cdot \frac{\mathbb{D}_M}{\mathbb{D}_B} \in F_{k-1, n-k}$$

Таким образом в качестве статистики критерия берем функцию

$$K = \frac{n-k}{k-1} \cdot \frac{\mathbb{D}_M}{\mathbb{D}_B}$$

В результате получили критерий: пусть $t_{кр}$ это квантиль распределения $F_{k-1, n-k}$ уровня значимости α , тогда

$$\begin{cases} H_0, & K < t_{кр} \\ H_1, & K \geq t_{кр} \end{cases}$$

1.9. Лекция 24.04.04.

Исследование статистической зависимости

Пусть случайная величина X зависит от величины Z (не обязательно случайной).

Def 1.9.1. Регрессией X на Z называется функция

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

Она показывает зависимость среднего значения X от значения Z . Уравнение $x = f(z)$ называется уравнением регрессии, а ее график — линей регрессии.

Пусть при n экспериментах при значениях z_1, \dots, z_n величины Z наблюдались соответствующие значения x_1, \dots, x_n случайной величины X . Обозначим через $\varepsilon_i = x_i - f(z_i)$ разницу между экспериментальными и теоретическими значениями X . Тогда $x_i = f(z_i) + \varepsilon_i$, где ε_i можно считать ошибкой наблюдения, случая и влияния неучтенный факторов.

Замечание 1.9.2. Обычно можно считать, что ε_i — независимые одинаковые нормальные случайные величины с нулевым первым моментом, т.к.

$$a = \mathbb{E}(\varepsilon_i) = \mathbb{E}(X_i) - \mathbb{E}(X | Z = z_i) = \mathbb{E}(X | Z = z_i) - \mathbb{E}(X | Z = z_i) = 0$$

Замечание 1.9.3. Вторым параметром σ^2 не всегда одинаковый, в некоторых ситуациях (временные ряды) ошибки ε_i могут быть зависимы.

Задача состоит в том, чтобы по данным значениям $(z_1, x_1), \dots, (z_n, x_n)$ как можно точнее оценить функцию регрессии $f(z)$. При этом предполагаем (часто из теории), что функция $f(z)$ определенного типа, но параметры которой не известны. В противном случае лучшим решением была бы любая кривая, проходящая через данные точки.

Метод наименьших квадратов

Этот метод состоит в выборе параметром функции $f(z)$ таким образом, чтобы сумма квадратов ошибок была наименьшей.

Def 1.9.4. Пусть $\Theta = (\Theta_1, \dots, \Theta_k)$ — набор неизвестных параметров функции $f(z)$. Оценка $\hat{\Theta}$, при которой достигается минимум $\sum \varepsilon_i^2$, называется оценкой метода наименьших квадратов.

Линейная парная регрессия

Пусть имеется линейная регрессия $f(z) = a + bz$, тогда $X_i = a + bz_i + \varepsilon_i$. Найдем оценки неизвестных параметров a и b методом наименьших квадратов.

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

Найдем частные производные.

$$\begin{aligned} \frac{\partial}{\partial a} \sum_{i=1}^n \varepsilon_i^2 &= 2 \sum_{i=1}^n (x_i - a - bz_i) \cdot (-1) = -2 \left(\sum_{i=1}^n x_i - \sum_{i=1}^n a - b \sum_{i=1}^n z_i \right) = -2(n\bar{x} - na - bn\bar{z}) \\ \frac{\partial}{\partial b} \sum_{i=1}^n \varepsilon_i^2 &= 2 \sum_{i=1}^n (x_i - a - bz_i) \cdot (-z_i) = -2 \left(\sum_{i=1}^n x_i z_i - \sum_{i=1}^n a z_i - b \sum_{i=1}^n z_i^2 \right) = -2(n\bar{x}\bar{z} - an\bar{z} - bn\bar{z}^2) \end{aligned}$$

Приравняем их к нулю.

$$\begin{cases} -2(n\bar{x} - na - bn\bar{z}) = 0 \\ -2(n\bar{x}\bar{z} - an\bar{z} - bn\bar{z}^2) = 0 \end{cases} \iff \begin{cases} a + b\bar{z} = \bar{x} \\ a\bar{z} + b\bar{z}^2 = \bar{x}\bar{z} \end{cases}$$

Получаем нормальную систему уравнений. Решим ее

$$\begin{cases} a = \bar{x} - b\bar{z} \\ (\bar{x} - b\bar{z})\bar{z} + b\bar{z}^2 = \bar{x}\bar{z} \end{cases} \iff \begin{cases} a = \bar{x} - b\bar{z} \\ b(\bar{z}^2 - (\bar{z})^2) = \bar{x}\bar{z} - \bar{x}\bar{z} \end{cases} \iff \begin{cases} \hat{a} = \bar{x} - b\bar{z} \\ \hat{b} = \frac{\bar{x}\bar{z} - \bar{x}\bar{z}}{\hat{\sigma}_z^2} \end{cases}$$

Получили оценки параметров a и b методов наименьших квадратов. Обозначим $\bar{x}_z = \mathbb{E}(X | Z = z)$ — условное среднее. Запишем уравнение линейной регрессии в более удобном виде

$$\begin{aligned} \bar{x}_z &= a + bz \\ \bar{x}_z &= \bar{x} - b\bar{z} + bz \\ \bar{x}_z - \bar{x} &= b(z - \bar{z}) \\ \bar{x}_z - \bar{x} &= \frac{\bar{x}\bar{z} - \bar{x}\bar{z}}{\hat{\sigma}_z^2} \cdot (z - \bar{z}) \\ \bar{x}_z - \bar{x} &= \frac{\bar{x}\bar{z} - \bar{x}\bar{z}}{\hat{\sigma}_x\hat{\sigma}_z} \cdot \frac{\hat{\sigma}_x}{\hat{\sigma}_z} \cdot (z - \bar{z}) \\ \bar{x}_z - \bar{x} &= \rho^* \cdot \frac{\hat{\sigma}_x}{\hat{\sigma}_z} \cdot (z - \bar{z}) \\ \frac{\bar{x}_z - \bar{x}}{\hat{\sigma}_x} &= \rho^* \cdot \frac{z - \bar{z}}{\hat{\sigma}_z} \end{aligned}$$

где

$$\rho^* = \frac{\bar{x}\bar{z} - \bar{x}\bar{z}}{\hat{\sigma}_x\hat{\sigma}_z}$$

это выборочный коэффициент линейной корреляции. Итого получили выборочное уравнение линейной регрессии

$$\frac{\bar{x}_z - \bar{x}}{\hat{\sigma}_x} = \rho^* \cdot \frac{z - \bar{z}}{\hat{\sigma}_z}$$

Замечание 1.9.5. При $n \rightarrow \infty$ имеем $\bar{x} \rightarrow \mathbb{E}(X)$, $\bar{z} \rightarrow \mathbb{E}(Z)$, $\bar{x}\bar{z} \rightarrow \mathbb{E}(xz)$ и следовательно выборочное уравнение регрессии стремится к

$$\frac{\mathbb{E}(X | Z = z) - \mathbb{E}(X)}{\sigma_x} = \rho \frac{z - \mathbb{E}(Z)}{\sigma_z}$$

Это теоретическое уравнение линейной регрессии, где

$$\rho = \frac{\mathbb{E}(xz) - \mathbb{E}(x)\mathbb{E}(z)}{\sigma_x\sigma_z}$$

это теоретический коэффициент линейной корреляции.

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

Выборочный коэффициент линейной корреляции

Def 1.9.7. Выборочным коэффициентом линейной корреляции называется величина

$$\rho^* = \frac{\bar{x}\bar{z} - \bar{x}\bar{z}}{\hat{\sigma}_x\hat{\sigma}_z}$$

Т.к. при замене в этой формуле средних на математические ожидания и выборочных средних квадратических отклонений на теоретические получаем теоретический коэффициент линейной корреляции ρ_{xz} , то данная величина является его точной оценкой. Таким образом выборочный коэффициент линейной корреляции характеризует силу линейной связи между случайными величинами. Его знак показывает, является ли она прямой или обратной.

Для оценки коэффициента линейной корреляции обычно используется шкала Чеддока (1.9.8).

Проверка гипотезы о значимости выборочного коэффициента корреляции

Пусть двумерная случайная величина $\langle z, x \rangle$ распределена нормально. По выборке объема n найден выборочный коэффициент линейной корреляции ρ^* . Теоретический коэффициент линейной корреляции обозначим ρ . Проверяется основная гипотеза H_0 о том, что $\rho = 0$, т.е. коэффициент ρ^* статистически не значим, против альтернативной гипотезы H_1 о том, что $\rho \neq 0$, т.е. коэффициент ρ^* статистически значим.

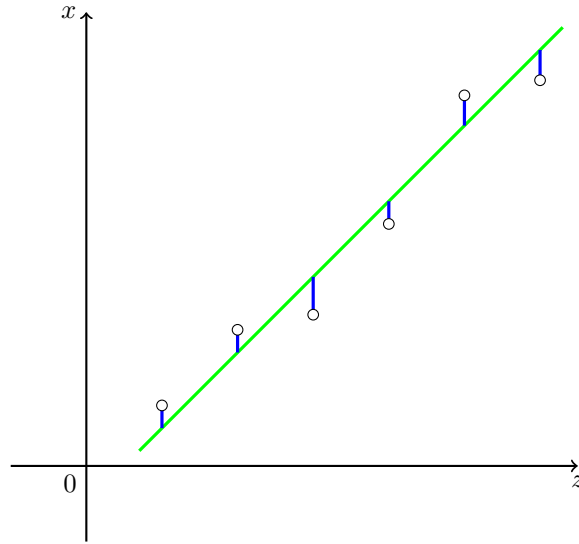


Рис. 1.9.6: Геометрический смысл прямой линейной регрессии

ρ^*	Характеристика силы связи
0.1—0.3	Слабая
0.3—0.5	Умеренная
0.5—0.7	Заметная
0.7—0.9	Сильная
> 0.9	Очень сильная

Таблица 1.9.8: Шкала Чеддока

Теорема 1.9.9. Если гипотеза H_0 верна, то

$$K = \frac{\rho^* \sqrt{n-2}}{\sqrt{1 - (\rho^*)^2}} \in T_{n-2}$$

На основе этой теоремы получаем критерий: пусть $t_{\text{кр}}$ это квантиль распределения $|T_{n-2}|$ уровня значимости α , тогда

$$\begin{cases} H_0, & |K| < t_{\text{кр}} \\ H_1, & |K| \geq t_{\text{кр}} \end{cases}$$

Исторический смысл понятия регрессия

Исследовалась зависимость роста детей от роста родителей (Гальтон, ≈ 1886 год). На основе собранных данных были получены следующие формулы

$$\begin{aligned} \mathbb{E}(P_s \mid z_f = z_1; z_m = z_2) &= 0.27z_1 + 0.2z_2 + \text{const} \\ \mathbb{E}(P_d \mid z_f = z_1; z_m = z_2) &= \frac{1}{1.08}P_s \end{aligned}$$

Казалось бы, что чем выше родители, тем выше должны быть их дети (прямая корреляция). В целом так и есть, но есть исключение: дети самых высоких родителей как правило были среднего роста, т.е. все сходилось к среднему значению. Отсюда и название «эффект регрессии» и «уравнение регрессии».

Выборочное корреляционное отношение

Пусть имеется k выборок случайной величины X при значениях z_1, \dots, z_k фактора Z .

Def 1.9.10. Выборочным корреляционным отношением называется величина

$$\eta_{xz} = \sqrt{\frac{\mathbb{D}_M}{\mathbb{D}_O}}$$

Lm 1.9.11.

$$0 \leq \eta_{xz} \leq 1$$

□ По 1.8.12

$$\left. \begin{array}{l} \mathbb{D}_O = \mathbb{D}_M + \mathbb{D}_B \\ \mathbb{D}_O \geq 0 \\ \mathbb{D}_B \geq 0 \\ \mathbb{D}_M \geq 0 \end{array} \right\} \Rightarrow 0 \leq \mathbb{D}_M \leq \mathbb{D}_O \Rightarrow 0 \leq \frac{\mathbb{D}_M}{\mathbb{D}_O} \leq 1$$

■

Lm 1.9.12. Если $\eta_{xz} = 1$, то имеется функциональная зависимость $x = f(z)$.

□ Если $\eta_{xz} = 1$, то $\mathbb{D}_M = \mathbb{D}_O$ и $\mathbb{D}_B = 0$. Т.е. при определенном значении z случайная величина X всегда принимает одно и то же значение, что является функциональной зависимостью. ■

Lm 1.9.13. Если $\eta_{xz} = 0$, то корреляция отсутствует.

□ Если $\eta_{xz} = 0$, то $\mathbb{D}_M = 0$, т.е. при разных уровнях z получили одно и то же выборочное среднее, что и означает отсутствие корреляции. ■

Lm 1.9.14. Имеет место неравенство $\eta \geq |\rho^*|$. Причем $\eta = |\rho^*|$ тогда и только тогда, когда имеет место точная корреляционная зависимость (все точки экспериментальных данных (z_i, x_i) лежат на одной прямой).

1.10. Лекция 24.04.11.

Ковариация и ее свойства

Def 1.10.1. Ковариацией случайных величин X и Y называется

$$\text{cov}(X, Y) = \mathbb{E}((X - \mathbb{E}(X))(Y - \mathbb{E}(Y))) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$$

Она является индикатором наличия и направления связи.

Пусть получены экспериментальные данные $(x_1, y_1), \dots, (x_n, y_n)$ случайных величин X и Y .

Def 1.10.2. Выборочной ковариацией называется величина

$$\text{cov}^*(X, Y) = \frac{1}{n} \sum_{i,j} (x_i - \bar{x})(y_j - \bar{y}) = \frac{1}{n} \sum_{i,j} x_i y_j - \bar{x} \bar{y} = \overline{xy} - \bar{x} \bar{y}$$

Замечание 1.10.3. Ясно, что выборочная ковариация является точной оценкой теоретической, но это смещенная оценка. Несмещенной оценкой будет $\frac{n}{n-1} \text{cov}^*(X, Y)$.

Свойства ковариации

1. $\text{cov}(X, Y) = \text{cov}(Y, X)$
2. $\text{cov}(X, a) = 0$, где $a = \text{const}$
3. $\text{cov}(X, bY) = b \text{cov}(X, Y)$
4. $\text{cov}(X + Y, Z) = \text{cov}(X, Z) + \text{cov}(Y, Z)$
5. $\text{cov}(X, X) = \mathbb{D}(X)$
6. $\mathbb{D}(X + Y) = \mathbb{D}(X) + \mathbb{D}(Y) + 2 \text{cov}(X, Y)$

Замечание 1.10.4. Эти же свойства имеет и выборочная ковариация.

Анализ модели парной линейной регрессии

Пусть при n экспериментах получены значения $(z_1, x_1), \dots, (z_n, x_n)$ величин X (случайная) и Z (фактор, не обязательно случайная). Пусть $X = \alpha + \beta Z + \varepsilon$ это теоретическая модель линейной регрессии. Случайный член ε отражает влияние неучтенных факторов, возможную нелинейность модели, ошибок измерения и просто случая. Цель заключается в том, чтобы дать оценки неизвестным параметрам α , β и ε .

Пусть при обработке данных методом наименьших квадратов нашли выборочное уравнение линейной регрессии $\hat{X} = a + b\hat{Z}$. Тогда $X_i = a + bZ_i + \varepsilon_i$, где $\varepsilon_i = X_i - a - bZ_i$ — наблюдаемые ошибки. a и b это точные оценки неизвестных параметров α и β .

Свойства ошибок ε_i

Lm 1.10.5.

$$\bar{\varepsilon}_i = 0$$

□

$$a = \bar{x} - b\bar{z} \implies a + b\bar{z} = \bar{x}$$

$$\bar{\varepsilon}_i = \overline{X_i - a - bZ_i} = \overline{X_i} - \overline{a + bZ_i} = \bar{x} - \bar{x} = 0$$

■

Lm 1.10.6.

$$\text{cov}(\hat{x}, \varepsilon) = 0$$

□

$$b = \frac{\overline{xz} - \bar{x}\bar{z}}{\mathbb{D}(Z)} = \frac{\text{cov}(X, Z)}{\mathbb{D}(Z)} \implies b\mathbb{D}(Z) - \text{cov}(X, Z) = 0$$

Значит

$$\text{cov}(\hat{x}, \varepsilon) = \text{cov}(a + bZ, X - a - bZ) = \text{cov}(bZ, X - bZ) = b\text{cov}(Z, X) - b^2\text{cov}(Z, Z) = b(\text{cov}(Z, X) - b\mathbb{D}(Z)) = 0$$

■

Анализ дисперсии результатовДисперсия наблюдаемых значений X_i :

$$\mathbb{D}(X) = \frac{1}{n} \sum_i (x_i - \bar{x})^2$$

Дисперсия расчетных значений $\hat{X}_i = a + bZ_i$:

$$\mathbb{D}(\hat{X}) = \frac{1}{n} \sum_i (\hat{x}_i - \bar{x})^2$$

Дисперсия экспериментальных ошибок:

$$\mathbb{D}(\varepsilon) = \frac{1}{n} \sum_i (\varepsilon_i - \bar{\varepsilon})^2 = \frac{1}{n} \sum_i \varepsilon_i^2$$

Знаем, что $X_i = a + bZ_i + \varepsilon_i$, т.е. $X_i = \hat{X}_i + \varepsilon_i$. Применим дисперсию к обеим частям, и получим

$$\mathbb{D}(X) = \mathbb{D}(\hat{X}) + \mathbb{D}(\varepsilon) + 2\text{cov}(\hat{X}, \varepsilon) = \mathbb{D}(\hat{X}) + \mathbb{D}(\varepsilon)$$

Получили разложение дисперсии. Ясно, что качество модели будет тем лучше, чем меньше дисперсия ошибок (в этом и суть метода наименьших квадратов). Значит модель будет тем лучше, чем больше доля первого слагаемого.

Def 1.10.7. Коэффициентом детерминации называется величина

$$R^2 = \frac{\mathbb{D}(\hat{X})}{\mathbb{D}(X)} = 1 - \frac{\mathbb{D}(\varepsilon)}{\mathbb{D}(X)}$$

Смысл коэффициента детерминации заключается в том, что это доля дисперсии объясненной при помощи данной модели. Величину $\frac{\mathbb{D}(\varepsilon)}{\mathbb{D}(X)}$ можно трактовать как долю необъясненной дисперсии.

Замечание 1.10.8. Ясно, что $0 \leq R^2 \leq 1$, причем чем больше R^2 , тем выше качество модели. Если $R^2 = 1$, то $\mathbb{D}(\varepsilon) = 0$ и, т.к. $\bar{\varepsilon} = 0$, то $\forall i \mid \varepsilon_i = 0$. Таким образом получили идеальную модель, все точки данных легли на прямую регрессии. Если $R^2 = 0$, то $\mathbb{D}(\hat{X}) = 0$ и $\hat{X} = \bar{X}$, т.е. построили ничего не объясняющую модель.

Проверка гипотезы о значимости уравнения регрессииПроверяется основная гипотеза H_0 о том, что $R^2_{\text{т}} = 0$ (т.е. R^2 статистически не значим), против альтернативной H_1 о том, что $R^2_{\text{т}} \neq 0$ (т.е. R^2 статистически значим).**Теорема 1.10.9.** Если нулевая гипотеза верна, то

$$K = \frac{R^2(n-2)}{1-R^2} \in F_{1,n-2}$$

На основе этой теоремы получаем критерий: пусть $t_{кр}$ это квантиль распределения $F_{1,n-2}$ уровня значимости α , тогда

$$\begin{cases} H_0, & |K| < t_{кр} \\ H_1, & |K| \geq t_{кр} \end{cases}$$

Замечание 1.10.10. Т.к. $R^2 = 0 \iff b = 0$, то это одновременно равносильно проверке гипотезы H_0 о том, что $\beta = 0$.

Связь между коэффициентом детерминации и коэффициентом линейной корреляции

Lm 1.10.11.

$$\sqrt{R^2} = \rho_{\hat{X}, X}$$

где $\rho_{\hat{X}, X}$ это коэффициент линейной корреляции между расчетными и наблюдаемыми значениями.

□

$$\text{cov}(\hat{X}, X) = \text{cov}(\hat{X}, \hat{X} + \varepsilon) = \text{cov}(\hat{X}, \hat{X}) + \underbrace{\text{cov}(\hat{X}, \varepsilon)}_{=0} = \mathbb{D}(\hat{X})$$

Таким образом

$$\rho_{\hat{X}, X} = \frac{\text{cov}(\hat{X}, X)}{\sqrt{\mathbb{D}(\hat{X}) \mathbb{D}(X)}} = \frac{\mathbb{D}(\hat{X})}{\sqrt{\mathbb{D}(\hat{X}) \mathbb{D}(X)}} = \sqrt{\frac{\mathbb{D}(\hat{X})}{\mathbb{D}(X)}} = \sqrt{R^2}$$

■

Lm 1.10.12.

$$\rho_{\hat{X}, X} = \rho_{X, Z}$$

□

$$\begin{aligned} \text{cov}(\hat{X}, X) &= \text{cov}(a + bZ, X) = b \text{cov}(Z, X) \\ \mathbb{D}(\hat{X}) &= \mathbb{D}(a + bZ) = b^2 \mathbb{D}(Z) \end{aligned}$$

Значит

$$\rho_{\hat{X}, X} = \frac{\text{cov}(\hat{X}, X)}{\sqrt{\mathbb{D}(\hat{X}) \mathbb{D}(X)}} = \frac{b \text{cov}(Z, X)}{\sqrt{b^2 \mathbb{D}(Z) \mathbb{D}(X)}} = \frac{\text{cov}(Z, X)}{\sqrt{\mathbb{D}(Z) \mathbb{D}(X)}} = \rho_{X, Z}$$

■

Замечание 1.10.13. В случае парной линейной регрессии коэффициент детерминации совпадает с квадратом коэффициента корреляции, т.е. $R^2 = \rho_{Z, X}^2$

Замечание 1.10.14. В случае парной линейной регрессии совпадают результаты проверок гипотез $H_0: R_t^2 = 0$, $H_0: \rho_{Z, X} = 0$ и $H_0: \beta = 0$.

Теорема Гаусса—Маркова

Пусть $X_i = \alpha + \beta Z_i + \varepsilon_i$ это теоретическая модель линейной регрессии, а $\hat{X} = a + bZ$ это выборочная уравнение линейной регрессии, полученное методом наименьших квадратов. Хотим узнать, насколько хороши оценки a и b неизвестных параметров α и β .

Теорема 1.10.15. (Гаусса—Маркова) Пусть выполнены следующие условия:

1. Случайные члены ε_i независимы и имеют одинаковое нормальное распределение $N(0; \sigma^2)$.
2. Случайные величины Z_i и ε_i независимы.

Тогда оценка (a, b) является наилучшей линейной несмещенной оценкой неизвестных параметров α и β , т.е.

1. Состоятельность: $a \xrightarrow{P} \alpha$ и $b \xrightarrow{P} \beta$ при $n \rightarrow \infty$.
2. Несмещенность: $\mathbb{E}(a) = \alpha$ и $\mathbb{E}(b) = \beta$.
3. Эффективность: a и b имеют наименьшую дисперсию в классе линейных оценок, равную

$$\mathbb{D}(a) = \frac{\overline{z^2} \sigma^2}{n \mathbb{D}(Z)} \quad \mathbb{D}(b) = \frac{\sigma^2}{n \mathbb{D}(Z)}$$

Замечание 1.10.16. Если ε_i зависимы или имеют разные дисперсии, то оценки по методу наименьших квадратов становятся неэффективными.

Замечание 1.10.17. Если случайные величины Z_i и ε_i зависимы, то оценки становятся смещенными и могут быть даже несостоятельными.

Стандартные ошибки коэффициентов регрессии

Видим, что дисперсии $\mathbb{D}(a)$ и $\mathbb{D}(b)$ зависят от дисперсии σ^2 случайного члена. По данным выборки экспериментальных ошибок ε_i получаем оценку дисперсии ошибок

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i \varepsilon_i^2$$

Однако это оценка смещенная, т.к.

$$\mathbb{E} \left(\frac{1}{n} \sum_i \varepsilon_i^2 \right) = \frac{n-2}{n} \sigma^2$$

поэтому несмещенной оценкой дисперсии σ^2 будет величина

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

Def 1.10.18. Величина S это стандартная ошибка регрессии. Она характеризует разброс результата вокруг линии регрессии.

Из этой формулы и 1.10.15 получем оценки дисперсий $\mathbb{D}(a)$ и $\mathbb{D}(b)$:

$$S_a^2 = \frac{\bar{z}^2 S^2}{n \mathbb{D}(z)} \quad S_b^2 = \frac{S^2}{n \mathbb{D}(z)}$$

Def 1.10.19. Коэффициенты S_a и S_b называются стандартными ошибками коэффициентов регрессии.

Прогнозирование в регрессионных моделях

Пусть $X_i = \alpha + \beta Z_i + \varepsilon_i$ это теоретическая модель, а $\hat{X} = a + bZ$ это модель метода наименьших квадратов. С помощью данной модели надо дать прогноз значения X_p при данном значении фактор Z_p . Тогда $X_p = \alpha + \beta Z_p + \varepsilon_i$ это реальное значение, а $\hat{X}_p = a + bZ_p$ — его точечная оценка (точечный прогноз). Обозначим $\Delta_p = \hat{X}_p - X_p$ ошибку прогноза.

Lm 1.10.20.

$$\mathbb{E}(\Delta_p) = 0$$

□ **TODO:** самостоятельно

■

Lm 1.10.21.

$$\mathbb{D}(\Delta_p) = \left(1 + \frac{1}{n} + \frac{(Z_p - \bar{z})^2}{n \mathbb{D}(Z)} \right) \cdot \sigma^2$$

где $\sigma^2 = \mathbb{D}(\varepsilon)$ это дисперсия случайного члена.

Замечание 1.10.22. Если σ^2 заменить на S^2 , то получаем стандартную ошибку прогноза

$$S_{\Delta_p} = S \sqrt{1 + \frac{1}{n} + \frac{(Z_p - \bar{z})^2}{n \mathbb{D}(Z)}}$$

Замечание 1.10.23. Проанализируем полученное выражение для $\mathbb{D}(\Delta_p)$.

1. Точность прогноза ограничена значением σ^2 , т.е. $\mathbb{D}(\Delta_p) \geq \sigma^2$.
2. $\mathbb{D}(\Delta_p) \rightarrow \sigma^2$ при $n \rightarrow \infty$, т.е. чем больше объем выборки, тем более качественная модель.
3. Чем дальше Z_p от \bar{Z} , тем хуже качество прогноза (рис. 1.10.24). Наилучшая точность достигается при $Z_p = \bar{Z}$, тогда $\mathbb{D}(\Delta_p) = \left(1 + \frac{1}{n} \right) \sigma^2$.

Доверительные интервалы прогноза и коэффициентов уравнения линейной регрессии

Пусть t_γ это квантиль двустороннего распределения Стьюдента с $n-2$ степенями свободы уровня γ . Тогда доверительные интервалы надежности γ для параметров α и β имеют вид

$$(a - t_\gamma S_a; a + t_\gamma S_a) \quad (b - t_\gamma S_b; b + t_\gamma S_b)$$

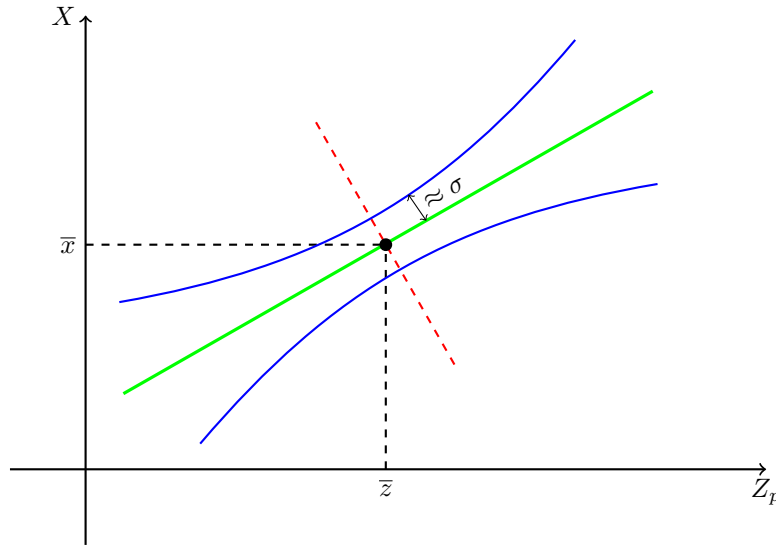


Рис. 1.10.24: Точность прогноза

Доверительный интервал прогноза надежности γ имеет вид

$$\left(\hat{X}_p - t_\gamma S_{\Delta p}; \hat{X}_p + t_\gamma S_{\Delta p} \right)$$

1.11. Лекция 24.04.18.

Математическое ожидание и дисперсия случайного вектора

Пусть имеется случайный вектор $\vec{X} = (X_1, \dots, X_n)^T$, где случайные величины X_i это компоненты случайного вектора.

Def 1.11.1. Математическим ожиданием случайного вектора \vec{X} называется вектор с координатами из математических ожиданий компонент.

$$\mathbb{E}(\vec{X}) = \begin{pmatrix} \mathbb{E}(X_1) \\ \vdots \\ \mathbb{E}(X_n) \end{pmatrix}$$

Def 1.11.2. Дисперсией (или матрицей ковариаций) случайного вектора \vec{X} называется матрица

$$\mathbb{D}(\vec{X}) = \mathbb{E}(\vec{X} - \mathbb{E}(\vec{X}))(\vec{X} - \mathbb{E}(\vec{X}))^T$$

Или, более просто, это матрица, состоящая из элементов $d_{i,j} = \text{cov}(X_i, X_j)$.

Замечание 1.11.3. На главной диагонали матрицы ковариаций стоят дисперсии компонент.

Свойства

1. $\mathbb{E}(A\vec{X}) = A\mathbb{E}(\vec{X})$, где A это матрица размера $n \times n$.
2. $\mathbb{E}(\vec{X} + \vec{B}) = \mathbb{E}(\vec{X}) + \vec{B}$
3. $\mathbb{D}(A\vec{X}) = A\mathbb{D}(\vec{X})A^T$
4. $\mathbb{D}(\vec{X} + \vec{B}) = \mathbb{D}(\vec{X})$

Общая модель линейной регрессии

Пусть результат X зависит от факторов Z_1, \dots, Z_k . Рассматривается теоретическая модель линейной регрессии.

$$\mathbb{E}(\vec{X} | \vec{Z}) = f(\vec{Z}) = \beta_1 Z_1 + \dots + \beta_k Z_k + \varepsilon$$

где $\vec{Z} = (z_1, \dots, z_k)^T$. Обозначим $\vec{\beta} = (\beta_1, \dots, \beta_k)^T$ — вектор неизвестных параметров регрессии, ε — случайный член (ошибка), отражающая нелинейность модели, влияние неучтенных факторов и т.д.

Пусть проведено $n \geq k$ экспериментов. Обозначим $\vec{Z}^{(i)} = (z_1^{(i)}, \dots, z_k^{(i)})$ набор значений факторов при i -ом эксперименте. Пусть $\vec{X} = (X_1, \dots, X_n)$ это набор результатов при n экспериментах. Согласно модели

$$\begin{cases} X_1 = \beta_1 z_1^{(1)} + \dots + \beta_k z_k^{(1)} + \varepsilon_1 \\ X_2 = \beta_1 z_1^{(2)} + \dots + \beta_k z_k^{(2)} + \varepsilon_2 \\ \vdots \\ X_n = \beta_1 z_1^{(n)} + \dots + \beta_k z_k^{(n)} + \varepsilon_n \end{cases}$$

где ε_i это теоретическая ошибка при i -ом эксперименте (она неизвестна). В матричной форме

$$\vec{X} = Z^T \vec{\beta} + \vec{\varepsilon}$$

где $\vec{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)$ это столбец ошибок, а матрица Z называется матрицей плана и имеет вид

$$Z_{k \times n} = \begin{pmatrix} z_1^{(1)} & z_1^{(2)} & \dots & z_1^{(n)} \\ z_2^{(1)} & z_2^{(2)} & \dots & z_2^{(n)} \\ \vdots & \vdots & \ddots & \vdots \\ z_k^{(1)} & z_k^{(2)} & \dots & z_k^{(n)} \end{pmatrix}$$

Требуется по данной матрице плана Z и вектору результатов \vec{X} найти оценки $\vec{B} = (b_1, \dots, b_k)$ для параметров регрессии $\vec{\beta} = (\beta_1, \dots, \beta_k)$ и параметров ошибок ε_i .

Замечание 1.11.4. В данной модели мы не теряем свободный член a , т.к. можно считать, что $z_1 \equiv 1$ и ей соответствует строка $(1, \dots, 1)$ в матрице плана Z .

Метод наименьших квадратов и нормальные уравнения

Будем считать, что выполнено условие о том, что $\text{rank } Z = k$, т.е. все строки матрицы плана независимы. Обозначим $A = ZZ^T$ — это будет квадратная матрица размера $k \times k$. Заметим, что

1. Матрица A симметричная, т.е. $A^T = A$.
2. Матрица A будет положительно определенной.
3. Существенная вещественная матрица \sqrt{A} такая, что $(\sqrt{A})^2 = A$.

Экспериментальная модель имеет вид

$$X_i = b_1 Z_{i1} + \dots + b_k Z_{ik} + \hat{\varepsilon}_i$$

где $\hat{\varepsilon}_i = X_i - (b_1 Z_{i1} + \dots + b_k Z_{ik})$.

При методе наименьших квадратов находим оценку $\vec{B} = (b_1, \dots, b_k)^T$, которая минимизирует функцию

$$L(\vec{B}) = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \|\hat{\varepsilon}\|^2 = \|\vec{X} - Z^T \vec{B}\|^2$$

Заметим, что $\|\vec{X} - Z^T \vec{B}\|^2$ это квадрат расстояния от точки $\vec{X} \in \mathbb{R}^n$ до точки $Z^T \vec{B} \in \mathbb{R}^n$. Причем $Z^T \vec{B}$ это произвольная точка подпространства $Z^T \vec{t}$, где $\vec{t} \in \mathbb{R}^k$. Заметим, что $\dim Z^T \vec{t} = k$, т.к. все строки матрицы плана независимы. Таким образом, искомое минимальное расстояние это расстояние от точки $\vec{X} \in \mathbb{R}^n$ до подпространства $Z^T \vec{t}$, $\vec{t} \in \mathbb{R}^k$. Это длина перпендикуляра, т.е. нужное расстояние получаем при условии, что вектор $\vec{X} - Z^T \vec{B}$ будет ортогонален всем векторам данного подпространства. Значит $\forall \vec{t} \in \mathbb{R}^k$ справедливо

$$(Z^T \vec{t}; \vec{X} - Z^T \vec{B}) = (Z^T \vec{t})^T \cdot (\vec{X} - Z^T \vec{B}) = \vec{t}^T Z (\vec{X} - Z^T \vec{B}) = \vec{t}^T (Z\vec{X} - ZZ^T \vec{B}) = \vec{t}^T (Z\vec{X} - A\vec{B}) = 0$$

Т.к. всем векторам данного пространства может быть ортогонален только нулевой вектор, то получаем, что

$$Z\vec{X} - A\vec{B} = \vec{0} \iff A\vec{B} = Z\vec{X} \iff \vec{B} = A^{-1}Z\vec{X}$$

Это нормальное уравнение или система нормальных уравнений с неизвестными b_1, \dots, b_k . Получили оценки по методу наименьших квадратов.

Свойства оценок по методу наименьших квадратов

Далее предполагаем, что выполнено условие $\text{rank } Z = k$ и ошибки ε_i — независимые нормальные случайные величины с распределением $N(0; \sigma^2)$. Таким образом $\mathbb{D}(\vec{\varepsilon}) = \sigma^2 E_n$.

Lm 1.11.5.

$$\vec{B} - \vec{\beta} = A^{-1}Z\vec{\varepsilon}$$

□

$$\begin{aligned}\vec{B} - \vec{\beta} &= A^{-1}Z\vec{X} - \vec{\beta} \\ &= A^{-1}Z\left(Z^T\vec{\beta} + \vec{\varepsilon}\right) - \vec{\beta} \\ &= A^{-1}ZZ^T\vec{\beta} + A^{-1}Z\vec{\varepsilon} - \vec{\beta} \\ &= A^{-1}A\vec{\beta} + A^{-1}Z\vec{\varepsilon} - \vec{\beta} \\ &= \vec{\beta} + A^{-1}Z\vec{\varepsilon} - \vec{\beta} \\ &= A^{-1}Z\vec{\varepsilon}\end{aligned}$$

Lm 1.11.6. Оценка \vec{B} это несмещенная оценка для параметра $\vec{\beta}$.

□

$$\mathbb{E}(\vec{B} - \vec{\beta}) = \mathbb{E}(A^{-1}Z\vec{\varepsilon}) = A^{-1}Z\mathbb{E}(\vec{\varepsilon}) = A^{-1}Z\vec{0} = \vec{0} \implies \mathbb{E}(\vec{B}) = \vec{\beta}$$

Lm 1.11.7.

$$\mathbb{D}(\vec{B}) = \sigma^2 A^{-1}$$

□

$$\begin{aligned}\mathbb{D}(\vec{B}) &= \mathbb{D}(\vec{B} - \vec{\beta}) \\ &= \mathbb{D}(A^{-1}Z\vec{\varepsilon}) \\ &= A^{-1}Z\mathbb{D}(\vec{\varepsilon})(A^{-1}Z)^T \\ &= A^{-1}Z\sigma^2 E_n Z^T A^{-1} \\ &= \sigma^2 A^{-1}ZZ^T A^{-1} \\ &= \sigma^2 A^{-1}AA^{-1} \\ &= \sigma^2 A^{-1}\end{aligned}$$

Замечание 1.11.8. Таким образом $\mathbb{D}(b_i) = \sigma^2 (A^{-1})_{i,i}$.

Введем обозначение $\hat{\sigma}^2$, которое определим как

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{1}{n} \|\vec{X} - Z^T \vec{B}\|^2 = \frac{1}{n} L(\vec{B})$$

Заметим, что $\hat{\sigma}^2$ это оценка неизвестной дисперсии ошибки σ^2 .

Теорема 1.11.9. Пусть выполнено условие $\text{rank } Z = k$ и ошибки ε_i — независимые нормальные случайные величины с распределением $N(0; \sigma^2)$. Тогда

1. Вектор $\frac{1}{\sigma} \sqrt{A}(\vec{B} - \vec{\beta})$ состоит из независимых случайных величин со стандартным нормальным распределением.

2.

$$\frac{n\hat{\sigma}^2}{\sigma^2} \in H_{n-k} \text{ и не зависит от } \vec{B}$$

3.

$$S^2 = \frac{n\hat{\sigma}^2}{(n-k)\sigma^2} = \frac{1}{n-k} \sum_{i=1}^n \hat{\varepsilon}_i^2 \quad \text{Это несмещенная оценка для } \sigma^2$$

□ **I пункт**

$$\frac{1}{\sigma} \sqrt{A}(\vec{B} - \vec{\beta}) = \frac{1}{\sigma} \sqrt{A}A^{-1}Z\vec{\varepsilon} = \frac{1}{\sigma} \sqrt{A^{-1}}Z\vec{\varepsilon}$$

Получили линейное преобразование нормального вектора $\vec{\varepsilon}$, и поэтому оно также является нормальным вектором.

$$\mathbb{E} \left(\frac{1}{\sigma} \sqrt{A} (\vec{B} - \vec{\beta}) \right) = \frac{1}{\sigma} \sqrt{A} \mathbb{E} (\vec{B} - \vec{\beta}) = \frac{1}{\sigma} \sqrt{A} \vec{0} = \vec{0}$$

Таким образом первый параметр компонент равен нулю.

$$\mathbb{D} \left(\frac{1}{\sigma} \sqrt{A} (\vec{B} - \vec{\beta}) \right) = \frac{1}{\sigma^2} \sqrt{A} \mathbb{D} (\vec{B} - \vec{\beta}) \sqrt{A^T} = \frac{1}{\sigma^2} \sqrt{A} \sigma^2 A^{-1} \sqrt{A} = E_k$$

Значит все координаты не коррелированы и имеют дисперсию равную единице.

II пункт

По построению метода наименьших квадратов имеем

$$\forall \vec{t} \in \mathbb{R}^k \mid \vec{X} - Z^T \vec{B} \perp Z^T \vec{t}$$

В частости это верно для $\vec{t} = \vec{B} - \vec{\beta}$. По обобщенной теореме Пифагора

$$\begin{aligned} \|\vec{X} - Z^T \vec{B}\|^2 + \|Z^T (\vec{B} - \vec{\beta})\|^2 &= \|\vec{X} - Z^T \vec{B} + Z^T (\vec{B} - \vec{\beta})\|^2 = \|\vec{X} - Z^T \vec{\beta}\|^2 \\ \|\vec{X} - Z^T \vec{B}\|^2 &= \|\vec{X} - Z^T \vec{\beta}\|^2 - \|Z^T (\vec{B} - \vec{\beta})\|^2 = \|\vec{\varepsilon}\|^2 - \|Z^T (\vec{B} - \vec{\beta})\|^2 \end{aligned} \quad (1)$$

Далее работает со вторым слагаемым.

$$\begin{aligned} \|Z^T (\vec{B} - \vec{\beta})\|^2 &= (\vec{B} - \vec{\beta})^T Z Z^T (\vec{B} - \vec{\beta}) \\ &= (\vec{B} - \vec{\beta})^T A (\vec{B} - \vec{\beta}) \\ &= (\vec{B} - \vec{\beta})^T \sqrt{A^T} \sqrt{A} (\vec{B} - \vec{\beta}) \\ &= \|\sqrt{A} (\vec{B} - \vec{\beta})\|^2 \\ &= \|\sqrt{A} A^{-1} Z \vec{\varepsilon}\|^2 \\ &= \|\sqrt{A^{-1}} Z \vec{\varepsilon}\|^2 \end{aligned} \quad (2)$$

Заметим, что строки матрицы $\sqrt{A^{-1}} Z$ ортогональны, т.к.

$$(\sqrt{A^{-1}} Z) (\sqrt{A^{-1}} Z)^T = \sqrt{A^{-1}} Z Z^T (\sqrt{A^{-1}})^T = \sqrt{A^{-1}} A \sqrt{A^{-1}} = E$$

$\sqrt{A^{-1}} Z$ это прямоугольная матрица размера $k \times n$. Из курса линейной алгебры известно, что ее можно дополнить до ортогональной матрицы C размера $n \times n$. Тогда первые k координат n -мерного вектора $\vec{Y} = \frac{1}{\sigma} C \vec{\varepsilon}$ совпадают с координатами вектора $\vec{Y} = \frac{1}{\sigma} \sqrt{A^{-1}} Z \vec{\varepsilon}$.

В результате из (1) и (2) получаем

$$\begin{aligned} \frac{n \hat{\sigma}^2}{\sigma^2} &= \frac{1}{\sigma^2} \sum_{i=1}^n \hat{\varepsilon}_i^2 \\ &= \frac{1}{\sigma^2} \|\vec{X} - Z^T \vec{B}\|^2 \\ &\stackrel{(1)}{=} \frac{1}{\sigma^2} \|\vec{\varepsilon}\|^2 - \frac{1}{\sigma^2} \|Z^T (\vec{B} - \vec{\beta})\|^2 \\ &= \left\| \frac{\vec{\varepsilon}}{\sigma} \right\|^2 - \left\| \frac{Z^T (\vec{B} - \vec{\beta})}{\sigma} \right\|^2 \\ &\stackrel{(2)}{=} \left\| \frac{\vec{\varepsilon}}{\sigma} \right\|^2 - \left\| \frac{\sqrt{A^{-1}} Z \vec{\varepsilon}}{\sigma} \right\|^2 \\ &= \sum_{i=1}^n \left(\frac{\varepsilon_i}{\sigma} \right)^2 - Y_1^2 - Y_2^2 - \dots - Y_k^2 \end{aligned} \quad (3)$$

Вектор $\frac{\vec{\varepsilon}}{\sigma}$ имеет n -мерное стандартное нормальное распределение, а Y_1, \dots, Y_k это первые k координат ортогонального преобразования данного вектора, поэтому (3) согласно 1.4.21 имеет распределение H_{n-k} и не зависит от вычитаемых координат вектора $\frac{\sqrt{A^{-1}} Z \vec{\varepsilon}}{\sigma}$, а значит и от вектора

$$\vec{B} = A^{-1}Z\vec{\varepsilon} + \vec{\beta} = \sigma\sqrt{A^{-1}}\left(\frac{1}{\sigma}\sqrt{A^{-1}}Z\vec{\varepsilon}\right)$$

которые являются их линейными комбинациями.

III пункт

Т.к. $\mathbb{E}(\chi_{n-k}^2) = n - k$, то

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

Это не равно σ^2 , т.е. оценка будет смещенной. Значит оценка

$$S^2 = \frac{n}{n - k} \hat{\sigma}^2 = \frac{1}{n - k} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

будет несмещенной. ■

1.12. Лекция 24.04.25.

Построение и анализ уравнения множественной линейной регрессии

Пусть выявлена зависимость результата X от факторов Z_1, \dots, Z_m . Проведено $n \geq m$ экспериментов и получены экспериментальные данные результата $\vec{X} = (X_1, \dots, X_n)$ при соответствующих значениях факторов $\vec{Z}^{(k)} = (Z_1^{(k)}, \dots, Z_m^{(k)})$. Предполагаем, что зависимость X от всех факторов линейная. Требуется по этим данным построить модель наилучшим образом объясняющую и предсказывающую поведение X .

Мультиколлинеарность

Def 1.12.1. Мультиколлинеарность это наличие заметной корреляции между всеми или некоторыми факторами.

Неприятные последствия мультиколлинеарности:

1. Оценки параметров становятся ненадежными, имеют большие стандартные ошибки и малую значимость, причем даже в том случае, когда модель в целом имеет высокую значимость.
2. Небольшое изменение исходных данных может привести к существенному изменению оценок регрессии.
3. Трудно выявить изолированное влияние конкретного фактора на результат и физический/экономический/т.д. смысл этого влияния.

Начальный отбор факторов в уравнение регрессии

Построим корреляционную матрицу, состоящую из коэффициентов линейной корреляции.

$$P = \begin{pmatrix} 1 & \rho_{x,z_1} & \rho_{x,z_2} & \dots & \rho_{x,z_m} \\ \rho_{z_1,x} & 1 & \rho_{z_1,z_2} & \dots & \rho_{z_1,z_m} \\ \rho_{z_2,x} & \rho_{z_2,z_1} & 1 & \dots & \rho_{z_2,z_m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_{z_m,x} & \rho_{z_m,z_1} & \rho_{z_m,z_2} & \dots & 1 \end{pmatrix}$$

Выбираем фактор, наиболее коррелирующий с X . Далее добавляем в модель факторы, которые с одной стороны имеют большую корреляцию с результатом X , а с другой стороны наименее коррелированы с факторами, которые уже включены в модель.

Пример 1.12.2. Пусть дана корреляционная таблица.

	X	Z_1	Z_2	Z_3
X	1	0.81	0.85	-0.65
Z_1	0.81	1	0.93	-0.38
Z_2	0.85	0.93	1	-0.28
Z_3	-0.65	-0.38	-0.28	1

Согласно алгоритму естественно включить фактор Z_2 , т.к. он имеет наибольшую корреляцию с X . Далее логично включить фактор Z_3 , т.к. Z_1 сильно коррелирован с уже включенным фактором Z_2 . Итого выбираем факторы Z_2 и Z_3 .

Анализ уравнения линейной регрессии

Пусть после отсева осталось k факторов. Теоретическая модель регрессии имеет вид

$$X = \beta_0 + \beta_1 z_1 + \dots + \beta_k z_k + \varepsilon$$

где ε — случайный член, отражающий влияние неучтенных факторов, возможную нелинейность модели, случай и т.д. Методом наименьших квадратов построили модель

$$\hat{X} = b_0 + b_1 z_1 + \dots + b_k z_k + \hat{\varepsilon}$$

Предполагаем, что $\forall i \mid \varepsilon_i \in N(0; \sigma^2)$ и независимы. Согласно пункту 3 в 1.11.9 получили несмещенную оценку для σ^2 в виде

$$S^2 = \frac{1}{n - k - 1} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

Def 1.12.3. S это стандартная ошибка регрессии.

Из 1.11.8 получили $\mathbb{D}(b_i) = \sigma^2 (A^{-1})_{i,i}$, где $A = ZZ^T$, а Z это матрица плана. Тогда

$$S_{b_i} = S \sqrt{(A^{-1})_{i,i}}$$

Это стандартные ошибки коэффициентов регрессии b_i .

Уравнение регрессии в стандартных масштабах

Замечание 1.12.4. При обычном уравнении регрессии по коэффициентам b_i нельзя судить о силе влияния фактора на результат X , т.к. факторы могут быть разной природы и иметь различные единицы измерения.

Стандартизация данных. Пусть имеется выборка (X_1, \dots, X_n) случайной величины X . Заменяем ее выборкой $t_x: X_i \rightarrow \frac{X_i - \bar{x}}{\hat{\sigma}_x}$, которую можно считать выборкой случайной величины $t_x = \frac{X - \mathbb{E}(X)}{\sigma_X}$, не имеющей единиц измерения.

Замечание 1.12.5. Очевидно, что $\bar{t}_x = 0$ и $\hat{\mathbb{D}}_{t_x} = 1$.

Lm 1.12.6.

$$\hat{\rho}_{X,Y} = \hat{\rho}_{t_x, t_y} = \overline{t_x t_y}$$

□

$$\hat{\rho}_{X,Y} = \frac{\text{cov}^*(X, Y)}{\hat{\sigma}_x \hat{\sigma}_y} = \frac{\frac{1}{n} - \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\hat{\sigma}_x \hat{\sigma}_y} = \frac{1}{n} \sum_{i=1}^n \frac{x_i - \bar{x}}{\hat{\sigma}_x} \frac{y_i - \bar{y}}{\hat{\sigma}_y} = \frac{1}{n} \sum_{i=1}^n t_x t_y = \overline{t_x t_y}$$

Равенство $\hat{\rho}_{t_x, t_y} = \overline{t_x t_y}$ будет выполняться, т.к. $\bar{t}_x = \bar{t}_y = 0$ и $\hat{\mathbb{D}}_{t_x} = \hat{\mathbb{D}}_{t_y} = 1$ в силу стандартизованности этих величин. ■

Все значения исследуемых признаков и результата стандартизируем

$$t_j = \frac{z_j^{(i)} - \bar{z}_j}{\hat{\sigma}_{z_j}} \quad 1 \leq i \leq k$$
$$t_x = \frac{x_i - \bar{x}}{\hat{\sigma}_x}$$

При линейной модели регрессии можно все величины в уравнении регрессии заменить на стандартизованные, т.к. все операции при стандартизации линейные. В результате получим так называемое уравнение в стандартных масштабах.

$$t_x = \gamma_1 t_1 + \dots + \gamma_k t_k$$

Замечание 1.12.7. Заметим, что свободного члена нет, т.к. $\bar{t}_x = \bar{t}_1 = \dots = \bar{t}_k = 0$.

Lm 1.12.8. При стандартизации система нормальных уравнений приобретает более простой вид.

$$\begin{cases} \gamma_1 + \rho_{z_1, z_2} \gamma_2 + \dots + \rho_{z_1, z_k} \gamma_k = \rho_{z_1, x} \\ \rho_{z_2, z_1} \gamma_1 + \gamma_2 + \dots + \rho_{z_2, z_k} \gamma_k = \rho_{z_2, x} \\ \dots \\ \rho_{z_k, z_1} \gamma_1 + \rho_{z_k, z_2} \gamma_2 + \dots + \gamma_k = \rho_{z_k, x} \end{cases}$$

В матричной форме это можно записать как $P\Gamma = P_x$, где P это матрица корреляций между факторами, $\Gamma = (\gamma_1, \dots, \gamma_k)^T$ и $P_x = (\rho_{x, z_1}, \dots, \rho_{x, z_k})^T$.

□ Действительно, нормальное уравнение регрессии имело вид $AB = Z\vec{X}$, где $A = ZZ^T$, а Z это матрица плана. Допустим, что все данные стандартизированны, тогда i -тый элемент столбца $Z\vec{X}$ имеет вид

$$\left(Z_1^{(i)}, \dots, Z_n^{(i)} \right) \cdot \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix} = n\overline{Z^{(i)}X} = n\rho_{z_i, x}$$

При этом элемент $a_{i,j}$ матрицы A будет равен

$$a_{i,j} = \left(Z_1^{(i)}, \dots, Z_n^{(i)} \right) \cdot \begin{pmatrix} Z_1^{(j)} \\ \vdots \\ Z_n^{(j)} \end{pmatrix} = \begin{cases} n\overline{z_i z_j} = n\rho_{z_i, z_j} & i \neq j \\ n\overline{z_i^2} = n\mathbb{D}(z_i) = n & i = j \end{cases}$$

После сокращения n в левой и правой части получаем искомую систему уравнений. ■

Переход от уравнения в стандартных масштабах к обычному уравнению регрессии и обратно можно выполнить по следующим формулам.

$$b_i = \gamma_i \frac{\sigma_x}{\sigma_{z_i}} \quad b_0 = \bar{X} - \sum_{i=1}^n b_i \bar{z_i}$$

Замечание 1.12.9. В случае парной линейной регрессии уравнение в стандартных масштабах имеет вид $t_x = \rho t_z$.

Смысл стандартизованных коэффициентов

Коэффициент γ_i на какую часть своего среднеквадратического отклонения σ_x изменится результат X при изменении фактора Z_i на величину своего среднеквадратического отклонения σ_{z_i} (при неизменных значениях других факторов). При мультиколлинеарности факторы оказывают не только прямое влияние на результат, но и косвенное (через влияние на другие факторы). Стандартизованный коэффициент γ_i можно трактовать как показатель прямого влияния, а остальные слагаемые в уравнении как результат косвенного влияния.

$$\rho_{z_i, z_1} \gamma_1 + \rho_{z_i, z_2} \gamma_2 + \dots + \rho_{z_i, z_{i-1}} \gamma_{i-1} + \gamma_i + \rho_{z_i, z_{i+1}} \gamma_{i+1} + \dots + \rho_{z_i, z_k} \gamma_k = \rho_{z_i, x}$$

Результат $\rho_{z_i, x}$ показывает величину полного влияния, а γ_i можно грубо трактовать как величину прямого влияния. Все остальное можно трактовать как косвенное влияние.

Замечание 1.12.10. Для измерения тесноты линейной связи между фактором и результатом при устранении влияния остальных факторов есть понятия коэффициентов частной корреляции. Например, при $k = 2$

$$\rho_{x, z_1 | z_2} = \frac{\rho_{x, z_1} - \rho_{x, z_2} \rho_{z_1, z_2}}{\sqrt{(1 - \rho_{z_1, z_2}^2) (1 - \rho_{x, z_2}^2)}} \quad \rho_{x, z_2 | z_1} = \frac{\rho_{x, z_2} - \rho_{x, z_1} \rho_{z_1, z_2}}{\sqrt{(1 - \rho_{z_1, z_2}^2) (1 - \rho_{x, z_1}^2)}}$$

Коэффициенты детерминации и множественной корреляции

Допустим, что как и в случае парной линейной регрессии дисперсию результата X можно разложить на объясненную и необъясненную составляющую, т.е.

$$\mathbb{D}(X) = \mathbb{D}(\hat{X}) + \mathbb{D}(\hat{\varepsilon})$$

где $\mathbb{D}(\hat{X})$ это дисперсия расчетных значений по построенной модели МНК, а $\mathbb{D}(\hat{\varepsilon})$ это дисперсия экспериментальных ошибок.

Def 1.12.11. Коэффициентом детерминации R^2 называется величина

$$R^2 = 1 - \frac{\mathbb{D}(\hat{\varepsilon})}{\mathbb{D}(X)}$$

Замечание 1.12.12. Ясно, что $0 \leq R^2 \leq 1$, причем чем больше R^2 , тем лучше качество модели. Если $R^2 = 1$, то $\mathbb{D}(\hat{\varepsilon}) = 0$, т.е. $\hat{\varepsilon} \equiv 0$, значит все экспериментальные данные легли на гиперплоскость регрессии. Если $R^2 = 0$, то $\mathbb{D}(\hat{X}) = 0$, т.е. $\hat{X} = \bar{X}$, значит $b_0 = \bar{X}, b_1 = \dots = b_k = 0$. Такая модель ничего не объясняет.

Замечание 1.12.13. В случае линейного уравнения регрессии $R^2 = \sum_{i=1}^k \gamma_i \rho_{x, z_i}$, где γ_i это стандартизованные коэффициенты.

Def 1.12.14. R называется коэффициентом множественной корреляции.

Замечание 1.12.15. При добавлении в модель новых факторов R^2 всегда вырастет, однако не всегда эти факторы следует вводить в модель, поэтому для выяснения того, следует ли это делать, существует скорректированный коэффициент детерминации \overline{R}^2 .

$$\overline{R}^2 = 1 - \frac{n-1}{n-k-1} \cdot \frac{\mathbb{D}(\hat{\varepsilon})}{\mathbb{D}(X)}$$

где n это число экспериментов, а k это число факторов в модели.

F-тест: проверка гипотезы о значимости уравнения регрессии в целом

Проверяется основная гипотеза H_0 о том, что $R_T^2 = 0$ (т.е. уравнение в целом не значимо), против альтернативной H_1 о том, что $R_T^2 \neq 0$.

Теорема 1.12.16. Если нулевая гипотеза верна, то

$$F = \frac{R^2}{1-R^2} \cdot \frac{n-k-1}{k} \in F_{k,n-k-1}$$

Пусть $t_{кр}$ это квантиль распределения $F_{k,n-k-1}$ уровня значимости α , тогда

$$\begin{cases} H_0, & F < t_{кр} \\ H_1, & F \geq t_{кр} \end{cases}$$

T-тест: проверка гипотезы о значимости отдельного коэффициента регрессии

Проверяется основная гипотеза H_0 о том, что $\beta_i = 0$, против альтернативной H_1 о том, что $\beta_i \neq 0$.

Теорема 1.12.17. Если нулевая гипотеза верна, то

$$T_i = \frac{b_i}{S_{b_i}} \in T_{n-k-1}$$

Пусть $t_{кр}$ это квантиль распределения $|T_{n-k-1}|$ уровня значимости α , тогда

$$\begin{cases} H_0, & |T_i| < t_{кр} \\ H_1, & |T_i| \geq t_{кр} \end{cases}$$

Замечание 1.12.18. Данный критерий применяется для отсева несущественных факторов.

Замечание 1.12.19. При мультиколлинеарности может оказаться так, что все коэффициенты по отдельности статистически не значимы, в то время как модель в целом имеет высокую значимость.

1.13. Лекция 24.05.02.

Нюансы регрессионного анализа

Пусть имеется линейное уравнение регрессии в матричной форме $\vec{X} = Z^T \vec{\beta} + \vec{\varepsilon}$, где $\vec{X} = (X_1, \dots, X_n)$ это столбец результатов n экспериментов, $Z_{k \times n}$ это матрица плана, $\vec{\beta} = (\beta_1, \dots, \beta_k)^T$ это столбец неизвестных (теоретических) коэффициентов регрессии и $\vec{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$ это вектор теоретических ошибок, возникающих из-за возможной нелинейности модели, невключенных факторов, случая и ошибок измерения.

Ранее предполагали, что строки матрицы плана Z независимы, т.е. $\text{rank } Z = k$ (на практике почти всегда выполнено), а также $\forall i \mid \varepsilon_i \in N(0; \sigma)$ и независимы. На практике часто бывают ситуации, когда ошибки коррелированы между собой и дисперсия зависит от эксперимента.

Взвешенный метод наименьших квадратов

Пусть ошибки ε_i некоррелированы, имеют нормальное распределение с первым параметром 0, но их дисперсия зависит от эксперимента, т.е.

$$\begin{aligned} \text{cov}(\varepsilon_i, \varepsilon_j) &= 0 & \bar{\varepsilon}_i &= 0 & \mathbb{D}(\varepsilon_i) &= \sigma^2 v_i \\ \mathbb{D}(\vec{\varepsilon}) &= \sigma^2 \text{diag}(v_1, \dots, v_n) = \sigma^2 V \end{aligned}$$

В этом случае оценки по методу наименьших квадратов неэффективны. Логично придать наблюдениям с меньшей дисперсией больший вес. Обозначим $w_i = \frac{1}{v_i}$ и каждое уравнение из системы $\vec{X} = Z^T \vec{\beta} + \vec{\varepsilon}$ умножим на $\sqrt{w_i}$. Тогда

$$\tilde{X}_i = \sqrt{w_i} X_i \quad \tilde{Z}_i^{(j)} = \sqrt{w_j} Z_i^{(j)} \quad \tilde{\varepsilon}_i = \sqrt{w_i} \varepsilon_i$$

При этом

$$\widetilde{\varepsilon}_i = 0 \quad \mathbb{D}(\widetilde{\varepsilon}_i) = w_i \mathbb{D}(\varepsilon_i) = \frac{1}{v_i} \cdot \sigma^2 v_i = \sigma^2$$

Таким образом $\mathbb{D}(\vec{\varepsilon}) = \sigma^2 E$, т.е. получили классическую ситуацию, и при применении к этим данным метода наименьших квадратов получаем эффективные оценки неизвестных коэффициентов регрессии. Рассмотрим несколько приложений этой модели.

I. Модель $\vec{X} = \beta_0 Z + \vec{\varepsilon}$

$$\vec{X} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \vec{\beta} = \beta_0 \quad \vec{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad Z^T = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix}$$

Подставим эти данные в нормальные уравнения $A\vec{B} = Z\vec{X}$, где

$$Z = (z_1, \dots, z_n) \quad A = ZZ^T = z_1^2 + \dots + z_n^2 \quad Z\vec{X} = (z_1, \dots, z_n) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = z_1 x_1 + \dots + z_n x_n$$

Таким образом нормальные уравнения приобретают вид

$$\sum_{i=1}^n z_i^2 \beta_0 = \sum_{i=1}^n z_i x_i$$

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n z_i x_i}{\sum_{i=1}^n z_i^2}$$

При взвешенном методе наименьших квадратов получаем оценку

$$\widetilde{\beta}_0 = \frac{\sum_{i=1}^n w_i z_i x_i}{\sum_{i=1}^n w_i z_i^2}$$

Пример 1.13.1. Рассмотрим модель $X = \beta_0 + \varepsilon$, т.е. $Z \equiv (1, \dots, 1)$. Тогда

$$\widetilde{\beta}_0 = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

Это можно трактовать как результат измерений скоропортящегося измерительного инструмента.

Пример 1.13.2. Пусть проводится n серий по k_i измерений (повторные измерения). Обозначим \bar{x}_i это среднее арифметическое измерений в i -ой серии, а ε_i — ошибка среднего измерения в i -ой серии. Тогда дисперсия $\mathbb{D}(\varepsilon_i) = \frac{\sigma^2}{k_i}$, где σ^2 это дисперсия ошибки при одном измерении. Тогда $w_i = k_i$ и получаем формулу

$$\widetilde{\beta}_0 = \frac{\sum k_i x_i}{\sum k_i}$$

Пример 1.13.3. Пусть X это потери тепла в квартире. Логично предположить, что основной фактор это разница температур снаружи и внутри. Рассмотрим уравнение $X = \beta Z + \varepsilon$, где Z это разница температур. Логично предположить, что с возрастанием Z возрастает дисперсия ошибки. Рассмотрим две ситуации.

I. $\mathbb{D}(\varepsilon_i) = cz_i$

$$\mathbb{D}(\varepsilon_i) = \sigma^2 \cdot \frac{cz_i}{\sigma^2} \implies w_i = \frac{\sigma^2}{cz_i} \implies \widetilde{\beta}_0 = \frac{\sum \frac{\sigma^2}{cz_i} z_i x_i}{\sum \frac{\sigma^2}{cz_i} z_i^2} = \frac{\sum x_i}{\sum z_i} = \frac{\bar{x}}{\bar{z}}$$

II. $\mathbb{D}(\varepsilon_i) = cz_i^2$

$$\mathbb{D}(\varepsilon_i) = \sigma^2 \cdot \frac{cz_i^2}{\sigma^2} \implies w_i = \frac{\sigma^2}{cz_i^2} \implies \widetilde{\beta}_0 = \frac{\sum \frac{\sigma^2}{cz_i^2} z_i x_i}{\sum \frac{\sigma^2}{cz_i^2} z_i^2} = \frac{\sum \frac{x_i}{z_i}}{\sum 1} = \overline{\left(\frac{x_i}{z_i}\right)}$$

Коррелированные наблюдения

Пусть дисперсии ошибок $\mathbb{D}(\varepsilon_i)$ различны и ошибки ε_i коррелированы между собой, причем эта корреляция известна $\text{cov}(\varepsilon_i, \varepsilon_j) = \sigma^2 v_{i,j}$. Тогда матрица ковариаций ошибок будет равна $\mathbb{D}(\vec{\varepsilon}) = \sigma^2 V$, где $V = (v_{i,j})$. Т.к. матрица V симметричная и положительно определенная, то существует \sqrt{V} . Умножим обе части модели на $(\sqrt{V})^{-1}$ слева. В результате получаем новую модель

$$\begin{aligned}\vec{X} &= \tilde{Z}^T \vec{\beta} + \vec{\tilde{\varepsilon}} \\ \vec{X} &= \sqrt{V^{-1}} \vec{X} \quad \tilde{Z}^T = \sqrt{V} Z^T \quad \vec{\tilde{\varepsilon}} = \sqrt{V} \vec{\varepsilon}\end{aligned}$$

Заметим, что $\tilde{\varepsilon}_i = 0$ и

$$\mathbb{D}(\vec{\tilde{\varepsilon}}) = \mathbb{D}(\sqrt{V^{-1}} \vec{\varepsilon}) = \sqrt{V^{-1}} \mathbb{D}(\varepsilon) (\sqrt{V^{-1}})^T = \sqrt{V^{-1}} \sigma^2 V \sqrt{V^{-1}} = \sigma^2 E$$

Таким образом получили стандартную ситуацию теоремы 1.10.15 и при этих данных получаем эффективные оценки неизвестных коэффициентов регрессии.

Нелинейная регрессия

Помимо общего метода наименьших квадратов многие нелинейные зависимости могут быть сведены к линейным при помощи простых приемов. Рассмотрим несколько случаев.

I. $X = \alpha + \beta f(Z) + \varepsilon$

Пусть $f(Z)$ это известная функция, например $f(Z) = \ln Z$ (быстрозамедляющийся процесс). Вычислим новые данные $Z'_i = f(Z_i)$ и получим стандартную модель парной линейной регрессии $X = \alpha + \beta Z' + \varepsilon$.

II. Степенная $X = \alpha Z^\beta + \varepsilon$

Прологарифмируем обе части и получим

$$\begin{aligned}\ln X &= \ln \alpha + \beta \ln Z + \ln \varepsilon \\ X' &= \alpha' + \beta Z' + \varepsilon'\end{aligned}$$

Далее по найденным α' и β' находим исходные α и β .

III. Показательная $X = \alpha e^{\beta Z} + \varepsilon$

Прологарифмируем обе части и получим

$$\begin{aligned}\ln X &= \ln \alpha + \beta Z + \ln \varepsilon \\ X' &= \alpha' + \beta Z + \varepsilon'\end{aligned}$$

IV. Полиномиальная $X = \alpha + \beta_1 Z + \dots + \beta_k Z^k + \varepsilon$

Введем новые переменные, положим $u_i = Z^i$, тогда

$$X = \alpha + \beta_1 u_1 + \dots + \beta_k u_k + \varepsilon$$

Полученное уравнение можно рассматривать как уравнение общей (множественной) регрессии.

Замечание 1.13.4. На практике обычно $k \leq 4$ во избежание мультиколлинеарности.

Замечание 1.13.5. При выборе между несколькими моделями выбираем ту, где коэффициент детерминации больше (т.е. меньше дисперсия экспериментальных ошибок).

Замечание 1.13.6. Построение даже удачной регрессионной модели не означает выявления причинно-следственной связи. Одна из причин заключается в том, что не учтен скрытый фактор.

Пример 1.13.7. Строилась модель точности бомбометания. Были выбраны факторы Z_1 — высота, Z_2 — ветер и Z_3 — количество истребителей противника. В итоге была получена модель

$$X = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2 + \beta_3 Z_3 + \varepsilon$$

где коэффициент β_3 оказался больше нуля, что противоречит логике. Это произошло из-за того, что в модель не был включен фактор Z_4 — облачность.

Замечание 1.13.8. Пусть эксперимент управляем и матрицу плана можно выбирать (почти) произвольно. Наиболее эффективные оценки получаются в том случае, когда строки матрицы плана ортогональны.

Замечание 1.13.9. Иногда вместо метода наименьших квадратов используется метод главных осей. Допустим, что есть корреляционное облако в виде эллипса. Согласно этому методу уравнение регрессии строим так, чтобы прямая совпадала с главной осью эллипса.

1.14. Лекция 24.05.16.

Датчики случайных чисел

Def 1.14.1. Члены последовательности y_1, \dots, y_n, \dots , которые можно рассматривать как экспериментальные значения данной случайной величины, называются псевдослучайными числами, а устройства или алгоритмы для их получения — датчиками случайных чисел.

I. Физические датчики

Теорема 1.14.2. Случайная величина $\eta \in U(0; 1)$ тогда и только тогда, когда разряды ξ_i ее двоичной записи $\eta = \sum_{i=1}^{\infty} 2^{-i} \xi_i$ — схема Бернулли с вероятностью $p = 0.5$.

Замечание 1.14.3. Недостатки физических датчиков:

1. Необходимо оборудование.
2. При повторном опыте нельзя получить ту же самую последовательность (а иногда это надо).

II. Таблицы случайных чисел

Пусть имеется таблица чисел (обычно двузначных) — результат работы некоторого датчика случайных чисел. Далее наугад выбирались строка и столбец и начиная с этого места брались «случайные» числа. При необходимости увеличения точности дописывали числа из соседнего столбца.

Замечание 1.14.4. Недостатки:

1. Требуется много оперативной памяти.
2. Слишком большая предсказуемость.

III. Математические датчики

Обычно это рекуррентные последовательности вида $y_n = f(y_{n-1})$. Основным считается мультипликативный датчик, который задается так: k_0 — начальное число, с которого начинается генерация, a — это множитель и m это модуль. Причем модуль выбирается так, что множитель a и k_0 взаимно просты с ним. В итоге последовательность псевдослучайных чисел строится по формулам

$$\begin{cases} k_n = k_{n-1} \cdot a \pmod{m} \\ y_n = \frac{k_n}{m} \in (0; 1) \end{cases}$$

Замечание 1.14.5. Рекомендации (старые, для 32-ух битных компьютеров):

1. $m = 2^{31} - 1$
2. $a = 630360016$ или $a = 764261123$
3. k_0 практически неважно.

IV. Датчик Уичмана и Хилла (1982)

Одновременно запускаются три мультипликативных датчика с параметрами $a_1 = 171$, $m_1 = 30269$, $a_2 = 172$, $m_2 = 30307$, $a_3 = 170$, $m_3 = 30323$. На i -ом шаге генерируется три псевдослучайных числа y'_n, y''_n, y'''_n . Тогда $y_n = \{y'_n + y''_n + y'''_n\}$ (дробная часть от суммы). Преимущества:

1. Он быстрее предыдущего.
2. Период этого датчика примерно $3 \cdot 10^{13}$, а предыдущего $2 \cdot 10^9$.

Моделирование случайных величин

В начале рассмотрим непрерывное распределение. Будем использовать метод обратной функции (квантильное преобразование).

Теорема 1.14.6. Пусть $F(x)$ — функция распределения абсолютно непрерывной случайной величины. Если $\eta \in U(0; 1)$, то случайная величина $\xi = F^{-1}(\eta)$ имеет функцию распределения $F(x)$.

Пример 1.14.7 (Показательное распределение). Пусть $\xi \in E_\alpha$, тогда

$$F(x) = \begin{cases} 0, & x < 0 \\ 1 - e^{-\alpha x}, & x \geq 0 \end{cases}$$

Значит получаем

$$y = 1 - e^{-\alpha x} \implies x = -\frac{1}{\alpha} \ln(1 - y) = F^{-1}(y)$$

Тогда согласно теореме, если y_i это значение датчика, то $x_i = F^{-1}(y) \in E_\alpha$.

Пример 1.14.8 (Нормальное распределение). Пусть $\xi \in N(a; \sigma^2)$. Известно, что если $\xi \in N(0; 1)$, то $(\sigma\xi + a) \in N(a; \sigma)$, поэтому достаточно уметь моделировать стандартное нормальное распределение.

$$F_0(x) = \frac{1}{2\sqrt{\pi}} \int_{-\infty}^x \exp\left(-\frac{z^2}{2}\right) dz$$

Таким образом $x_i = F_0^{-1}(y_i) \in N(0; 1)$, если $y_i \in U(0; 1)$.

Замечание 1.14.9. Преимущество этого метода заключается в простоте и универсальности, а недостаток заключается в не очень высокой эффективности (раньше было так, сейчас уже возможно иначе).

Нормальные случайные числа

I. На основе центральной предельной теоремы

Пусть $\eta_i \in U(0; 1)$, тогда $\mathbb{E}(\eta_i) = a = 0.5$ и $\mathbb{D}(\eta_i) = \frac{1}{12}$. По центральной предельной теореме

$$\frac{S_n - na}{\sqrt{n\mathbb{D}(\xi)}} = \frac{S_n - \frac{n}{2}}{\sqrt{\frac{n}{12}}} \Rightarrow N(0; 1)$$

Уже при $n = 12$ получаются довольно хорошие результаты.

II. Точное моделирование пары случайных величин с распределением $N(0; 1)$

Теорема 1.14.10. Пусть $\eta_1, \eta_2 \in N(0; 1)$ и независимы. Тогда следующие величины $X, Y \in N(0; 1)$ и независимы.

$$X = \sqrt{-2 \ln \eta_1} \cos(2\pi\eta_2) \quad Y = \sqrt{-2 \ln \eta_1} \sin(2\pi\eta_2)$$

□ Пусть $X, Y \in N(0; 1)$ и независимы. Тогда плотность совместного распределения

$$f_{X,Y}(x, y) = \frac{1}{2\sqrt{\pi}} \exp\left(-\frac{x^2}{2}\right) \frac{1}{2\sqrt{\pi}} \exp\left(-\frac{y^2}{2}\right) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x^2 + y^2)\right)$$

При переходе к полярным координатам плотность данного распределения примет вид

$$f_{X,Y}(\rho, \varphi) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}\rho^2\right) \rho = \underbrace{\frac{1}{2\pi}}_{f_\varphi} \cdot \underbrace{\exp\left(-\frac{\rho^2}{2}\right) \rho}_{f_\rho}$$

Найдем их функции распределения.

$$F_\varphi(x) = \frac{x}{2\pi}$$

$$F_\rho(x) = \int_0^\rho \rho \exp\left(-\frac{\rho^2}{2}\right) d\rho = -\exp\left(-\frac{\rho^2}{2}\right) \Big|_0^\rho = 1 - \exp\left(-\frac{\rho^2}{2}\right)$$

Отсюда методом обратной функции получаем формулы для моделирования случайных величин φ и ρ .

$$y = \frac{x}{2\pi} \implies x = 2\pi y \quad y \in U(0; 1)$$

$$y = 1 - \exp\left(-\frac{\rho^2}{2}\right) \implies \rho = \sqrt{-2 \ln(1 - y)} \quad y \in U(0; 1)$$

Т.к. $y \in U(0; 1)$, то $1 - y \in U(0; 1)$. В итоге получаем

$$X = \rho \cos \varphi = \sqrt{-2 \ln \eta_1} \cos(2\pi\eta_2) \quad Y = \rho \sin \varphi = \sqrt{-2 \ln \eta_1} \sin(2\pi\eta_2)$$

■

Быстрый показательный датчик

Теорема 1.14.11. Пусть случайные величины $\eta_1, \dots, \eta_{2n-1} \in U(0; 1)$. Обозначим ξ_1, \dots, ξ_{n-1} это расставленные в порядке возрастания значения $\eta_{n+1}, \dots, \eta_{2n-1}$. Отдельно положим $\xi_0 = 0$ и $\xi_n = 1$. Тогда μ_i независимы и имеют показательное распределение с параметром α .

$$\mu_i = -\frac{1}{\alpha} (\xi_i - \xi_{i-1}) \ln(\eta_1 \cdot \dots \cdot \eta_n) \quad 1 \leq i \leq n$$

Пример 1.14.12. При $n = 3$ получаем $\eta_1, \eta_2, \eta_3, \eta_4, \eta_5$. Пусть $\eta_4 < \eta_5$ (в противном случае переобозначим их), тогда

$$\begin{aligned}\mu_1 &= -\frac{1}{\alpha} \eta_4 \ln(\eta_1 \eta_2 \eta_3) \\ \mu_2 &= -\frac{1}{\alpha} (\eta_5 - \eta_4) \ln(\eta_1 \eta_2 \eta_3) \\ \mu_3 &= -\frac{1}{\alpha} (1 - \eta_5) \ln(\eta_1 \eta_2 \eta_3)\end{aligned}$$

Замечание 1.14.13. При этом методе экономим время при вычислении логарифма, но теряем на сортировке. Алгоритм оптимален при $n = 3$. В этом случае он примерно в два раза быстрее метода обратной функции.

Моделирование дискретных распределений

Пусть имеется дискретное распределение, которое задается парами x_k и $p_k = P(\xi = c_k)$, где $k = 1, 2, \dots$. Разобьем отрезок $[0; 1]$ на отрезки длины p_k . Обозначим $\eta_m = \sum_{k=1}^m p_k$ и $\eta_0 = 0$ — границы отрезков. Пусть $y_i \in [0; 1]$ это случайное число. Если $y_i \in [\rho_{j-1}; \rho_j)$, то полагаем, что $x_i = c_j$.

I. Распределение Бернулли

Пусть $\xi \in B_p$, тогда получаем два отрезка, т.е.

$$\begin{cases} y_i \in [0; 1 - p) \implies x_i = 0 \\ y_i \in [1 - p; 1] \implies x_i = 1 \end{cases}$$

II. Биномиальное распределение

Пусть $\xi \in B_{n,p}$, т.е. $P(\xi = k) = C_n^k p^k q^{1-k}$, где $k = 0, \dots, n$. Исходя из смысла биномиального распределения $\xi = \xi_1 + \dots + \xi_n$, где $\xi_i \in B_p$. Далее берем n значений датчика $y_i \in U(0; 1)$ и

$$\begin{cases} y_i \in [0; 1 - p) \implies z_i = 0 \\ y_i \in [1 - p; 1] \implies z_i = 1 \end{cases} \implies x_k = \sum_{i=1}^n z_i$$

III. Геометрическое распределение

Пусть $\xi \in G_p$, т.е. $P(\xi = k) = q^{k-1}p$, где $k = 1, 2, \dots$. Исходя из смысла геометрического распределения ξ это номер первого успешного испытания, поэтому берем значения датчика $y_i \in U(0; 1)$ до тех пор, пока оно не попадет в интервал $[1 - p; 1]$. Далее полагаем $x_k = i$.

IV. Распределение Пуассона

Пусть $\xi \in \Pi_\lambda$, т.е. $P(\xi = k) = \frac{\lambda^k}{k!} e^{-\lambda}$, где $k = 0, 1, \dots$

Теорема 1.14.14. Пусть μ_1, μ_2, \dots — независимые случайные величины, имеющие показательное распределение с параметром λ . Положим $S_n = \mu_1 + \dots + \mu_n$ и $N = \max n \mid S_n \in [0; 1]$. Тогда $N \in \Pi_\lambda$.

На основе данной теоремы и метода обратной функции получаем формулу для моделирования распределение Пуассона.

$$n_i = \min \left\{ k \mid \prod_{j=1}^k y_{i,j} < e^{-\lambda} \right\}$$

где $y_{i,j}$ это значение датчика в i -ой серии.

1.15. Лекция 24.05.23.

Метод Монте—Карло

Цель метода состоит в том, чтобы находить неизвестные значения изучаемой величины при помощи «разыгрывания» некоторой случайной величины.

Общая постановка метода. Пусть требуется найти неизвестное число a и имеется случайная величина ξ такая, что $E(\xi) = a$. Тогда согласно закону больших чисел

$$\frac{\xi_1 + \dots + \xi_n}{n} \xrightarrow{\text{п.н.}} a$$

Следовательно, при достаточно больших n среднее выборочное $\bar{x} \approx a$.

Оценка погрешности. Пусть $D(\xi) < \infty$, тогда согласно центральной предельной теореме

$$\frac{S_n - na}{\sqrt{n\mathbb{D}(\xi_1)}} = \frac{n(\bar{x} - a)}{\sqrt{n\mathbb{D}(\xi_1)}} \Rightarrow z \in N(0; 1)$$

По правилу трех сигм $P(|z| < 3) \approx 0.9973$. Следовательно, при больших n можно считать, что

$$\left| \frac{n(\bar{x} - a)}{\sqrt{n\mathbb{D}(\xi_1)}} \right| < 3 \Rightarrow |(\bar{x} - a)| < 3 \frac{\sqrt{\mathbb{D}(\xi_1)}}{\sqrt{n}}$$

Отсюда видим, что сходимость есть, но она достаточно медленная, порядка $\frac{1}{\sqrt{n}}$, поэтому на практике не удастся получить очень точные оценки для a .

Вычисление определенных интегралов

$$\int_a^b \varphi(x) dx = \lim_{\Delta x_i \rightarrow 0} \sum_{i=1}^n \varphi(c_i) \Delta x_i$$

где Δx_i это длины интервалов разбиения, а c_i — точки внутри интервалов. На этом определении основаны так называемые квадратурные формулы.

I. Формула прямоугольников

Разобьем отрезок $[a; b]$ на n равных частей длины $\Delta x_i = \frac{b-a}{n}$. Обозначим x_i — середину i -ого интервала. Тогда

$$I = \int_a^b \varphi(x) dx \approx \frac{b-a}{n} \sum_{i=1}^n \varphi(x_i) = I_n$$

Можно показать, что $|I - I_n| \leq \frac{M_1}{n^2}$, где $M_1 = \text{const}$.

II. Формула трапеций

Разобьем отрезок $[a; b]$ на n равных частей $[x_i; x_{i+1}]$ длины $\Delta x_i = \frac{b-a}{n}$. Обозначим $y_i = \varphi(x_i)$, где $0 \leq i \leq n$. Тогда

$$I = \int_a^b \varphi(x) dx \approx \frac{b-a}{2n} (y_0 + y_n + 2(y_1 + \dots + y_{n-1})) = I_n$$

Можно показать, что $|I - I_n| \leq \frac{M_2}{n^2}$, где $M_2 = \text{const}$, но M_2 меньше, чем M_1 , примерно в два раза.

III. Формула Симпсона (формула парабол)

Разобьем отрезок $[a; b]$ на $n = 2m$ равных отрезков $[x_i; x_{i+1}]$ длины $\Delta x_i = \frac{b-a}{n}$. Обозначим $y_i = \varphi(x_i)$, где $0 \leq i \leq n$. Тогда

$$I = \int_a^b \varphi(x) dx \approx \frac{b-a}{3n} (y_0 + y_n + 4(y_1 + y_3 + \dots + y_{n-1}) + 2(y_2 + y_4 + \dots + y_{n-2})) = I_n$$

Можно показать, что $|I - I_n| \leq \frac{M}{n^4}$, где $M = \text{const}$.

IV. Метод Монте—Карло

В качестве узлов берутся псевдослучайные числа. Пусть $I = \int_0^1 \varphi(x) dx$. Ясно, что $[0; 1] \rightarrow [a; b]$ при помощи линейной замены. Обозначим $\eta_i \in U(0; 1)$ — значение датчика случайных чисел, $f_{\eta_i}(x) \equiv 1, x \in [0; 1]$. Пусть $\xi_i = \varphi(\eta_i)$, тогда

$$\mathbb{E}(\xi_i) = \int_{-\infty}^{\infty} \varphi(x) \cdot f_{\eta_i}(x) dx = \int_0^1 \varphi(x) dx = I$$

По методу Монте—Карло получаем, что

$$I \approx \hat{I}_n = \frac{1}{n} \sum_{i=1}^n \varphi(\eta_i)$$

где η_i это значения датчика случайных чисел. Погрешность вычислений не будет превосходить

$$|I - I_n| \leq \frac{3\sqrt{\mathbb{D}(\xi_1)}}{\sqrt{n}} \quad \mathbb{D}(\xi_1) = \int_0^1 \varphi^2(x) dx - I^2$$

Замечание 1.15.1. Скорость сходимости намного хуже, чем в квадратурных формулах. При этом для оценки погрешности потребуются вычислить (или оценить сверху) дополнительный интеграл, поэтому метод Монте—Карло не используется для вычисления определенных интегралов.

Кратные интегралы

Замечание 1.15.2. При вычислении k -кратных интегралов число узлов сетки возрастает как n^k и аналог метода прямоугольников будет довольно трудоемким. Таким образом метод Монте—Карло становится уместным, т.к. он не зависит от размерности.

Для вычисления интеграла

$$I = \int_0^1 \cdots \int_0^1 \varphi(x_1, \dots, x_k) dx_1 \dots dx_k$$

достаточно случайно бросить n точек в данный k -мерный единичный куб. Для этого последовательные k значений датчика берем в качестве координат соответствующей случайной точки. Вычисляем значение подынтегральной функции в этих случайных точках и получаем

$$I \approx I_n = \frac{1}{n} \sum_{i=1}^n \varphi(x_{i_1}, \dots, x_{i_k})$$

При этом скорость сходимости остается равной $\frac{1}{\sqrt{n}}$. В частности, если функция $\varphi = I_D$ это индикатор области D , то с помощью метода Монте—Карло можно приблизительно найти объем области D .

Пример 1.15.3. Пусть необходимо оценить площадь четверти единичной окружности. Генерируем последовательность псевдослучайных чисел и разбиваем их на пары $\langle x_{i_1}, x_{i_2} \rangle$, которые будут координатами точек.

$$I_D = \begin{cases} 0, & x_{i_1}^2 + x_{i_2}^2 > 1 \\ 1, & x_{i_1}^2 + x_{i_2}^2 \leq 1 \end{cases}$$

Итого $S = \frac{n_D}{n}$, где n_D это число точек, попавших в область D .

Метод расслоенной выборки

Пусть имеется k -мерный интеграл

$$I = \int_0^1 \cdots \int_0^1 \varphi(x_1, \dots, x_k) dx_1 \dots dx_k$$

Каждую из сторон k -мерного куба разобьем N равных частей, тогда куб разобьется на $n = N^k$ маленьких кубиков Δ_i со стороной $\frac{1}{N}$. В каждом из этих кубиков возьмем случайную точку, построенную с помощью датчика случайных чисел, $\eta_i = (\eta_{i_1}^1, \dots, \eta_{i_k}^k) \in \Delta_i$, где $1 \leq i \leq n$. Интеграл оценивается при помощи суммы

$$I_n = \frac{1}{n} \sum_{i=1}^n \varphi(\eta_i)$$

При этом методе погрешность будет лучше, а именно

$$|I - I_n| \leq C \cdot \frac{1}{n^{\frac{1}{2} + \frac{1}{k}}} \quad C = \text{const}$$

Замечание 1.15.4. В самых лучших квадратурных формулах для k -кратных интегралов погрешность составляет

$$|I - \hat{I}| \leq C \cdot \frac{1}{n^{1+\varepsilon}} \quad C = \text{const}$$

Равномерность по Вейлю

Def 1.15.5. Числовая последовательность x_1, \dots, x_n , где $x_i \in [0; 1]$ называется равномерной по Вейлю, если частота попадания точек x_i на любой отрезок $[a; b]$ стремится к его длине $b - a$ при $n \rightarrow \infty$.

Замечание 1.15.6. Ясно, что значения равномерного стандартного распределения или датчика случайных чисел будут равномерными по Вейлю.

Теорема 1.15.7. Пусть α это иррациональное число, тогда последовательность $x_n = \{n\alpha\}$ является равномерной по Вейлю.

Замечание 1.15.8. Если последовательность x_n равномерная по Вейлю, то

$$\frac{1}{n} \sum_{i=1}^n \varphi(x_i) \rightarrow \int_0^1 \varphi(x) dx$$

НО ВОЗМОЖНО ЧТО

$$\frac{1}{n} \sum_{i=1}^n \psi(x_{i_1}, x_{i_2}) \rightarrow C \neq \int_0^1 \int_0^1 \psi(x, y) dx dy$$

Def 1.15.9. Числовая последовательность x_1, \dots, x_n, \dots называется вполне равномерной по Вейлю, если для произвольного числа k частота попадания k -мерных точек $\langle x_i^{(1)}, \dots, x_i^{(k)} \rangle$ в любой параллелепипед внутри единичного куба стремится к его объему.

Замечание 1.15.10. Хороший датчик случайных чисел должен выдавать вполне равномерную последовательность. Мультипликативный датчик не является вполне равномерным.

Пример 1.15.11 (Парадокс первой цифры). Найти вероятность того, что число 2^n начинается с цифры 7.

Решение: Пусть 2^n начинается с цифры m . Это означает, что

$$\begin{aligned} m \cdot 10^l &\leq 2^n < (m+1) \cdot 10^l \\ \log_{10} m + l &\leq n \log_{10} 2 < \log_{10}(m+1) + l \in [l; l+1) \\ \log_{10} m &\leq \{n \log_{10} 2\} < \log_{10}(m+1) \end{aligned}$$

Согласно теореме 1.15.7 последовательность $\{n \log_{10} 2\}$ является равномерной по Вейлю, следовательно вероятность того, что число начнется с цифры m стремится к длине интервала, т.е. к $\log_{10}(m+1) - \log_{10} m = \log_{10}(1 + \frac{1}{m})$. Итого при $m = 7$ получаем $\log_{10}(1 + \frac{1}{7}) \approx 0.058$.

1.16. Лекция 24.05.30.

Энтропия

Пусть случайная величина ξ это результат некоторого эксперимента с N исходами A_1, \dots, A_N . Вероятности этих исходов — p_1, \dots, p_N .

Def 1.16.1. Энтропией эксперимента называется величина

$$H(\xi) = - \sum_{i=1}^N p_i \log_2 p_i$$

При $p_i = 0$ соответствующее слагаемое полагаем равным нулю.

Замечание 1.16.2. $H(\xi) \geq 0$, т.к. $\log_2 p_i < 0$, то каждое слагаемое будет неположительным. Минус перед знаком суммы дает общую неотрицательность энтропии.

Замечание 1.16.3. $H(\xi) = 0$ тогда и только тогда, когда $\exists i$ такое, что $p_i = 1$ и $p_k = 0$ при $k \neq i$. В этом случае результат эксперимента не случаен, а полностью предопределен (нет неопределенности результата).

Lm 1.16.4. Максимум энтропии, равный $\log_2 N$, достигается при $p_1 = \dots = p_N = \frac{1}{N}$.

□ Рассмотрим дискретную случайную величину η , значениями которой будут p_1, \dots, p_N , а вероятности этих значений все будут равны $\frac{1}{N}$. Рассмотрим функцию $\varphi(x) = x \log_2 x$ и найдем ее вторую производную.

$$\varphi''(x) = \left(\log_2 x + x \cdot \frac{1}{x \ln 2} \right)' = \frac{1}{x \ln 2}$$

Это больше нуля при $x > 0$, значит эта функция выпукла вниз. Тогда по неравенству Йенсена получаем

$$\begin{aligned} \varphi(\mathbb{E}(\eta)) &\leq \mathbb{E}(\varphi(\eta)) \\ \varphi(\mathbb{E}(\eta)) &= \varphi\left(\sum_{i=1}^N p_i \frac{1}{N}\right) = \varphi\left(\frac{1}{N}\right) = -\frac{1}{N} \log_2 N \\ \mathbb{E}(\varphi(\eta)) &= \sum_{i=1}^N \frac{1}{N} p_i \log_2 p_i = \frac{1}{N} \sum_{i=1}^N p_i \log_2 p_i = \frac{1}{N} H(\xi) \\ -\frac{1}{N} \log_2 N &\leq \frac{1}{N} H(\xi) \\ \log_2 N &\geq H(\xi) \end{aligned}$$

■

Замечание 1.16.5. Энтропию можно рассматривать как количественную характеристику неопределенности эксперимента. Если $H = 0$, то результат предопределен, а если $H = H_{\max} = \log_2 N$, то все исходы равновероятны и ни одному из них нельзя отдать предпочтение.

Пример 1.16.6. Пусть $\xi \in B_p$, тогда

$$H(\xi) = -(1-p) \log_2(1-p) - p \log_2 p$$

При $p = 0.5$ имеем

$$H(\xi) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

Пример 1.16.7. Перед испытуемым зажигалось n лампочек, и он должен был быстро указать на загоревшуюся лампочку. Оказалось, что самое длинное время реакции было в том случае, когда лампочки загорались с одинаковой частотой. Если частоты были разными, то время реакции было прямопропорционально энтропии эксперимента.

Кодирование информации

Пусть требуется передать сообщение из N символов. Пусть длины кодовых слов равны, тогда для каждого символа понадобится $\log_2 N$ бит, а для всего сообщения потребуется $N \log_2 N$ бит. Для больших по объему сообщений можно сократить число бит, учитывая что разные символы возникают с разной частотой. Пусть p_1, \dots, p_N — вероятности соответствующих N символов, тогда при длинном сообщении i -ый символ встречается $v_i = Np_i$ раз.

Def 1.16.8. Назовем сообщение типичным, если $\forall i: |v_i - Np_i| < \delta$.

Обозначим $M_{N,\delta}$ — число типичных сообщений.

Теорема 1.16.9. (Макмиллана (частный случай))

$$\frac{1}{N} \log_2 M_{N,\delta} \xrightarrow{n \rightarrow \infty} H$$

Отсюда число типичных сообщений $M_{N,\delta} \leq 2^{N(H+\varepsilon)}$, где $\varepsilon > 0$ — мало. Обозначим $H_0 = \log_2 N$, тогда с вероятностью, близкой к единице, можно сократить длину сообщения с коэффициентом сжатия $\gamma = \frac{H}{H_0}$. Если символы появляются независимо, то большее сжатие невозможно, а если использовать их зависимость, то коэффициент сжатия можно существенно уменьшить. Пусть γ_∞ — данный коэффициент сжатия. Для русского языка $\gamma \approx 0.87$, а $\gamma_\infty \approx 0.24$ (для литературного языка) и $\gamma_\infty \approx 0.17$ (для деловых сообщений).

Def 1.16.10. Величина $1 - \gamma_\infty$ называется избыточностью языка.

Энтропия абсолютно непрерывных распределений

Def 1.16.11. Пусть ξ это абсолютно непрерывная случайная величина с плотностью $f(x)$. Энтропией случайной величины ξ называется величина

$$H(\xi) = - \int_{-\infty}^{\infty} f(x) \log_2 f(x) dx$$

Теорема 1.16.12. Следующие распределения имеют наибольшую энтропию:

1. Если $\xi \in [0; 1]$, то $\xi \in U(0; 1)$.
2. Если $\xi \in [0; +\infty)$ и $\mathbb{E}(\xi) = 1$, то $\xi \in E_1$.
3. Если $\xi \in (-\infty; +\infty)$, $\mathbb{E}(\xi) = 0$ и $\mathbb{D}(\xi) = 1$, то $\xi \in N(0; 1)$.

Задача о разорении игрока

Играют два игрока. Вероятность выигрыша первого игрока — p , а второго — $1 - p$. Ставка в каждой игре это одна единица. Капитал первого игрока равен k , а второго — $M - k$. Игра прекращается, когда один из игроков потеряет свой капитал. Найти вероятность разорения первого игрока.

Выигрыш первого игрока $S_n = z_1 + \dots + z_n$, где $P(z_i = 1) = p$ и $P(z_i = -1) = q$. Пусть r_k это вероятность разорения при стартовом капитале k , тогда $r_k = pr_{k+1} + qr_{k-1}$. Это линейное однородное рекуррентное уравнение. Сначала рассмотрим случай $p \neq \frac{1}{2}$. Составим и решим характеристическое уравнение.

$$px^2 - x + (1-p) = 0 \implies x_{1,2} = \frac{1 \pm \sqrt{1-4p(1-p)}}{2p} = \frac{1 \pm (1-2p)}{2p} = \begin{cases} \frac{1-p}{p} = \lambda \\ 1 \end{cases} \implies r_k = c_1 1^k + c_2 \lambda^k$$

Это общее решение. Чтобы найти константы подставим начальные (в данном случае их лучше назвать граничными) условия.

$$\begin{cases} r_0 = 1 = c_1 + c_2 \\ r_M = 0 = c_1 + c_2 \lambda^M \end{cases} \iff \begin{cases} c_1 = 1 - c_2 \\ 1 - c_2 + c_2 \lambda^M = 0 \end{cases} \iff \begin{cases} c_2 = \frac{1}{1-\lambda^M} \\ c_1 = \frac{-\lambda^M}{1-\lambda^M} \end{cases}$$

Таким образом получаем итоговое решение

$$r_k = \frac{-\lambda^M}{1-\lambda^M} + \frac{\lambda^k}{1-\lambda^M} = \frac{\lambda^k - \lambda^M}{1-\lambda^M}$$

Далее рассмотрим ситуацию $M \rightarrow \infty$. Возможны два случая.

1. $p < \frac{1}{2}$, тогда $\lambda > 1$ и $\lambda^M \rightarrow \infty$, значит

$$r_k = \frac{\lambda^k - \lambda^M}{1 - \lambda^M} = \frac{\lambda^{k-M} - 1}{\lambda^{-M} - 1} = 1$$

Таким образом игрок гарантированно разоряется.

2. $p > \frac{1}{2}$, тогда $\lambda < 1$ и $\lambda^M \rightarrow 0$, значит

$$r_k = \frac{\lambda^k - \lambda^M}{1 - \lambda^M} = \lambda^k = \left(\frac{q}{p}\right)^k$$

Вернемся к решению характеристического уравнения и рассмотрим случай $p = \frac{1}{2}$. Тогда $x_1 = x_2 = 1$ и общее решение будет иметь вид

$$r_k = c_1 1^k + c_2 k \cdot 1^k = c_1 + c_2 k$$

Для нахождения констант подставим начальные условия.

$$\begin{cases} r_0 = 1 = c_1 \\ r_M = 0 = c_1 + c_2 M \end{cases} \iff \begin{cases} c_1 = 1 \\ c_2 = -\frac{1}{M} \end{cases} \implies r_k = 1 - \frac{k}{M}$$

При $M \rightarrow \infty$ получаем, что $r_k \rightarrow 1$. Значит вне зависимости от стартового капитала при бесконечной игре первый игрок рано или поздно разорится.