

Содержание

Лекция 1.	4
Выборки	4
Выборочные характеристики	4
Начальная обработка статданных	5
Геометрическая интерпретация данных	6
Лекция 2.	8
Точечная оценка	8
Свойство точечных оценок	8
Точечные оценки моментов	8
Метод моментов (Пирсона)	10
Лекция 3.	11
Метод максимального правдоподобия	11
Неравенство Рао-Крамера	13
Лекция 4.	15
Основные распределения математической статистики	15
Распределение «хи-квадрат»	15
Распределение Стюдента	16
Распределение Фишера-Снедекера	16
Математическое ожидание и дисперсия случайного вектора	16
Многомерное нормальное распределение	17
Многомерная центральная предельная теорема	18
Лемма Фишера	18
Основная теорема	18
Лекция 5.	20
Квантильное распределение	20
Интервальные оценки	20
Доверительные интервалы для параметров нормального распределения	21
Асимптотические доверительные интервалы	23
Лекция 6.	24
Проверка статистических гипотез	24
Построение критериев согласия	25
Гипотеза о среднем нормальной совокупности при известной дисперсии	25

Гипотеза о среднем нормальной совокупности при неизвестной дисперсии	26
Доверительные интервалы как критерии гипотез по параметрам распределения	26
Критерий вероятности появления события	27
Лекция 7.	27
Критерии для проверки гипотез о распределении	27
Простая параметрическая гипотеза	27
Сложная параметрическая гипотеза	28
Критерии для проверки однородности	29
Проверки однородности выборок из нормальных совокупностей	30
Лекция 8.	32
Статистическая зависимость	32
Корреляционное облако	32
Корреляционная таблица	32
Критерий «хи-квадрат» для проверки независимости	33
Однофакторный дисперсионный анализ	34
Общая, внутригрупповая и межгрупповая дисперсии	34
Проверка гипотезы о влиянии фактора	35
Лекция 9.	35
Исследование статистической корреляции	35
Математическая модель регрессии	35
Метод наименьших квадратов	36
Линейная парная регрессия	36
Геометрический смысл линии регрессии	38
Выборочный коэффициент линейной корреляции	38
Проверка гипотезы о значимости выборочного коэффициента корреляции	39
Выборочное корреляционное отношение	39
Лекция 10.	40
Свойство ковариации	40
Анализ модели линейной парной регрессии	40
Анализ дисперсии результата	41
Проверка гипотезы о значимости уравнения регрессии	42
Связь между коэффициентом детерминации и коэффициентом линейной корреляции	42
Теорема Гаусса-Маркова	43
Стандартные ошибки коэффициентов регрессии	43
Прогнозирование регрессионных моделей	44

Доверительные интервалы прогноза и коэффициентов уравнения линейной регрессии 44

Лекция 11. 45

Математическое ожидание и дисперсия случайного вектора 45

Уравнение общей регрессии 45

Метод наименьших квадратов и нормальные уравнения 46

Свойства оценок метода наименьших квадратов 47

Оценка дисперсии случайного члена 47

Лекция 12. 48

Построение и анализ уравнения множественной линейной регрессии 48

Мультиколлинеарность 48

Отбор факторов в уравнении регрессии 48

Анализ уравнения линейной регрессии 49

Уравнение регрессии стандартных масштабов 49

Смысл стандартизованных коэффициентов 51

Коэффициенты детерминации и множественной корреляции 51

Скорректированный коэффициент детерминации 52

Проверка гипотез по значимости уравнения регрессии 52

Лекция 13. 53

Невадия регрессионного анализа 53

Взвешенный метод наименьших квадратов 53

Коррелированные наблюдения 54

Составление матрицы плана при управляемом эксперименте 54

Метод главных осей 55

Нелинейные регрессии 55

Х. Программа экзамена в 2024/2025 57

Лекция 1.

Теория вероятности изучает характеристику случайных величин, тогда как математическая статистика решает обратную задачу

Допустим, что у нас есть случайная величина, по ней мы можем найти математическое ожидание, моменты и оценить, какое распределение имеет случайная величина.

Выборки

Def. Выборка - набор данных, полученных в ходе экспериментов. Тогда количество экспериментов n - объем Выборки

Def. Генеральной совокупностью называются все результаты проведенных экспериментов

Def. Выборочной совокупностью называются наблюдаемые данные экспериментов

Не все данные экспериментов мы можем наблюдать, например, выборы, тогда опросы голосовавших - выборочная совокупность, а результаты выборов - генеральная. Очевидно, что выборочная и генеральная совокупности могут иметь различные распределения.

Def. Выборка называется **репрезентативной**, если ее распределение близко к распределению генеральной совокупностью

Пример - **ошибка выжившего**. Во время Второй Мировой стал вопрос, в каких местах стоит бронировать корпус самолета. Самолеты возвращались с пулевыми отверстиями, и интуитивно казалось, что стоит бронировать те места, которые больше всего пострадали. Однако не были учтены те самолеты, которые не вернулись, а те, которые выжили, выжили благодаря тому, что были прострелены в нелетальных местах, поэтому было принято решение бронировать фюзеляж в менее пострадавших местах

В дальнейшем считаем, что все выборки репрезентативны

Def. 1. Выборкой объема n называется набор из n экспериментальных данных $\vec{X} = (x_1, x_2, \dots, x_n)$ (апостериорное определение)

Def. 2. Выборкой объема n называется набор из n независимых одинаково распределенных случайных величин $\vec{X} = (X_1, X_2, \dots, X_n)$ (априорное определение)

Выборочные характеристики

Можно выборку рассматривать как дискретную случайную величину с одинаковыми вероятностями $p_i = \frac{1}{n}$ и вычислить для нее математическое ожидание, дисперсию и функцию распределения

Def. Выборочным средним \bar{x} называется величина $\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i$

Def. Выборочной дисперсией D^* называется величина $D^* = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^2$ (или $D^* = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{x}^2$)

По закону больших чисел выборочное среднее будет сходиться к матожиданию

Def. Исправленной дисперсией называется величина $S^2 = \frac{n}{n-1} D^* = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{x})^2$

Def. Выборочной функцией распределения $F^*(x)$ называется функция $F^*(x) = \frac{\text{число данных } x_i < x}{n}$

Th. Выборочная функция распределения поточечно сходится к теоретической функции распределения:

$$\forall y \in \mathbb{R} F^*(y) \xrightarrow{p} F(y)$$

$$F(y) = P(X < y)$$

$$F_y^* = \frac{1}{n} \sum_{i=1}^n I(X_i < y) \xrightarrow[\text{по ЗБЧ}]{p} EI(X_i < y) = P(X_i < y) = P(X_1 < y) = F_{X_1}(y)$$

Усилим теорему

Th. Гливенко-Кантелли. $\sup_{x \in \mathbb{R}} |F^*(x) - F(x)| \xrightarrow{p} 0$

Th. Колмогорова. $\sqrt{n} \sup_{x \in \mathbb{R}} |F^*(x) - F(x)| \rightrightarrows K$ - распределение Колмогорова с функцией распределения $F_K(x) = \sum_{j=-\infty}^{\infty} (-1)^j e^{-2j^2 x^2}$, $x \in [0; \infty)$

Начальная обработка статданных

1. Ранжирование данных - упорядочиваем выборки по возрастанию. В результате получаем вариационный ряд $\vec{X} = (X_{(1)}, X_{(2)}, \dots, X_{(n)})$

$$X_{(1)} = \min X_i; \quad X_{(n)} = \max X_i$$

$X_{(i)}$ - i -ая порядковая статистика

2. Объединим повторяющиеся данные - получаем т.н. частотный вариационный ряд

X_i	$X_{(1)}$	\dots	$X_{(r)}$	\sum
n_i	n_1	\dots	n_r	n

Иногда часть данных отбрасывается сверху и снизу (по 5, по 10, по 5% и так далее), чтобы сделать выборку репрезентативной

$$\text{Тогда } \bar{x} = \frac{1}{n} \sum X_i n_i, \quad D^* = \frac{1}{n} \sum (X_i - \bar{x})^2 n_i$$

3. Чтобы уменьшить количество вычислений или сделать гистограмму, делают интервальный вариационный ряд: разбиваем данные на интервалы и считаем, сколько данных n_i попало в интервал.

Тогда n_i - частота интервала A_i

Есть два основных способа разбиения на интервалы:

- (а) Интервалы одинаковой длины
- (б) Равнонаполненные интервалы (в каждом интервале примерно одинаковое количество данных)

Число интервалов K такое, что $\frac{K(n)}{n} \rightarrow 0$ и $K(n) \xrightarrow{n \rightarrow \infty} \infty$

Обычно применяют формулу Стерджесса $K \approx 1 + \log_2 n$ или $K \approx \sqrt[3]{n}$

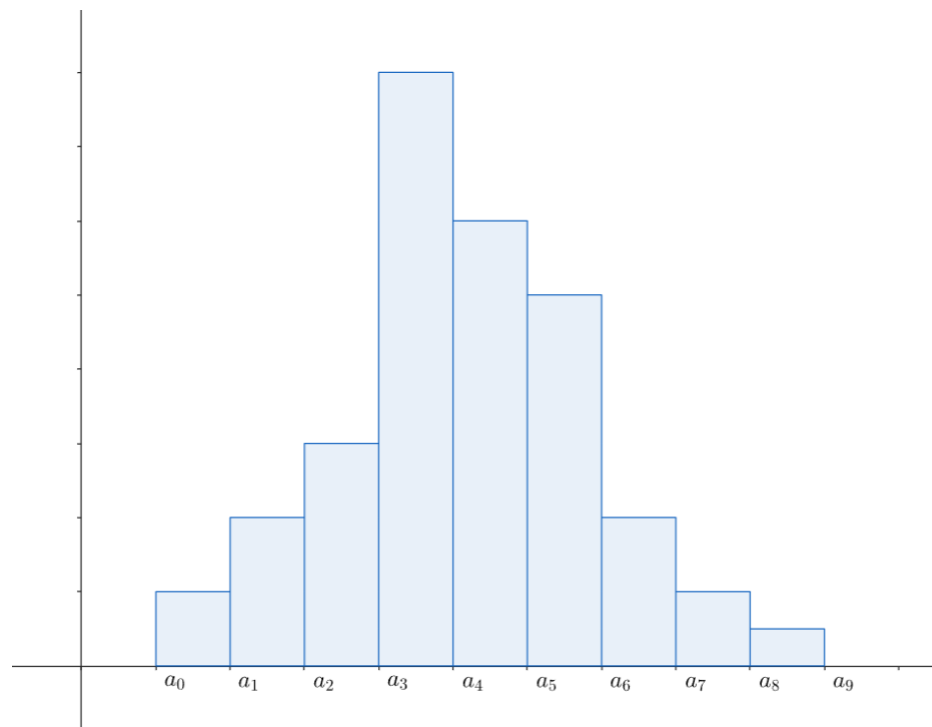
Пусть получили интервальный вариационный ряд

интервалы	$[a_0; a_1)$	$[a_1; a_2)$	\dots	$[a_{K-1}; a_K]$	\sum
частоты	n_1	n_2	\dots	n_K	n

Геометрическая интерпретация данных

- Гистограмма

Строится ступенчатая фигура из прямоугольников, основание i -ого прямоугольника - интервал, высота прямоугольника - $\frac{n_i}{nl_i}$, где l_i - длина интервала

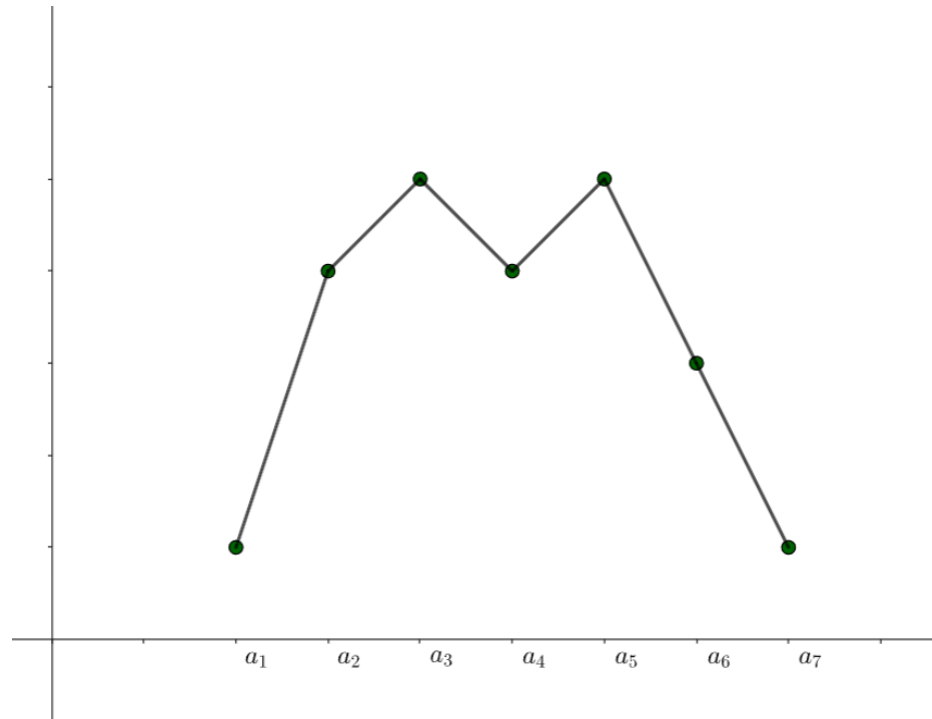


Визуально можно сделать гипотезу, как ведет себя распределение.

Th. Гистограмма поточечно сходится к теоретической плотности

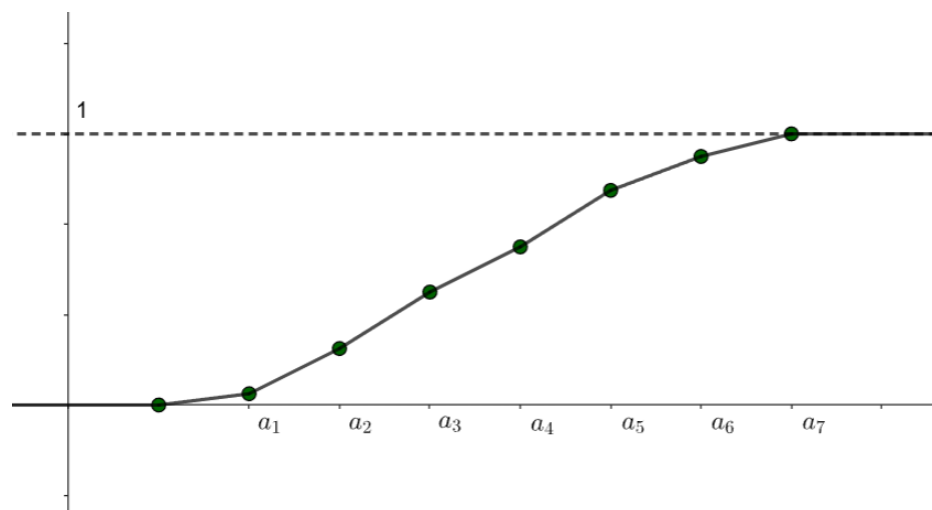
• Полигон

На оси абсцисс отмечаем значения частотного вариационного ряда, по оси ординат - их частоты. Получившиеся точки соединяем отрезками



• Выборочная функция распределения

На основе таблицы строится график функции распределения



Она может быть ступенчатой, ломаной или соединена по усмотрению

Лекция 2.

Точечная оценка

Пусть имеется выборка $\vec{X} = (X_1, X_2, \dots, X_n)$ объемом n

Пусть требуется найти приближенную оценку θ^* неизвестного параметра θ

Находим ее при помощи некоторой функции обработки данных $\theta^* = \theta^*(X_1, \dots, X_n)$

Def. Такая функция называется статистикой

Def. А оценка θ^* называется точечной оценкой

Свойство точечных оценок

1. Состоятельность

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

2. Несмещенность

Def. Оценка θ^* параметра θ называется несмещенной, если математическое ожидание $E\theta^* = \theta$

Nota. Оценка θ^* называется асимптотически несмещенной, если $E\theta^* \xrightarrow{p} \theta$ при $n \rightarrow \infty$

3. Эффективность

Def. Оценка θ_1^* не хуже θ_2^* , если $E(\theta_1^* - \theta)^2 \leq E(\theta_2^* - \theta)^2$. Или, если θ_1^* и θ_2^* несмещенные, то $D\theta_1^* \leq D\theta_2^*$

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

Nota. Не существует эффективной оценки в классе всех возможных оценок

Th. В классе несмещенных оценок существует эффективная оценка

4. Асимптотическая нормальность

Def. Оценка θ^* параметра θ называется асимптотически нормальной, если $\sqrt{n}(\theta^* - \theta) \Rightarrow N(0, \sigma^2(\theta))$ при $n \rightarrow \infty$

Точечные оценки моментов

Def. Выборочным средним \bar{x} называется величина $\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i$

Def. Выборочной дисперсией D^* называется величина $D^* = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^2$

Def. Исправленной дисперсией S^2 называется величина $S^2 = \frac{n}{n-1} D^* = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{x})^2$

Def. Выборочным средним квадратическим отклонением называется величина $\sigma^* = \sqrt{D^*}$

Def. Исправленным средним квадратическим отклонением называется величина $S = \sqrt{S^2}$

Def. Выборочным k -ым моментом называется величина $\bar{x}^k = \frac{1}{n} \sum_{i=1}^n X_i^k$

Def. Модой Mo^* называется варианта x_k с наибольшей частотой $n_k = \max_i (n_1, n_2, \dots, n_m)$

Def. Выборочной медианой Me^* называется варианта x_i в середине вариационного ряда

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

Th. \bar{x} - состоятельная несмещенная оценка теоретического матожидания $\hat{A}X = a$

1) $E\bar{x} = a$

2) $\bar{x} \xrightarrow{p} a$ при $n \rightarrow \infty$

1) $E\bar{x} = E\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{1}{n} \sum_{i=1}^n EX_i = \frac{1}{n} n EX_1 = EX_1 = a$

2) $\bar{x} = \frac{\bar{x}_1 + \dots + \bar{x}_n}{n} \xrightarrow{p} a$ согласно Закону Больших Чисел

Nota. Если второй момент конечен, то \bar{x} - асимптотически нормальная оценка. По ЦПТ $\frac{S_n - nEX_1}{\sqrt{n} \sqrt{DX_1}} = \sqrt{n} \frac{\bar{x} - EX_1}{\sqrt{DX_1}} \Rightarrow N(0, 1)$ или $\sqrt{n}(\bar{x} - EX_1) \Rightarrow N(0; DX_1)$

Th. Выборочный k -ый момент является состоятельной несмещенной оценкой теоретического k -ого момента

1) $\overline{EX^k} = EX^k$

2) $\overline{X^k} \xrightarrow{p} X^k$

Это следует из предыдущей теоремы, если взять X^k вместо X

Th. Выборочной дисперсией D^* и S^2 являются состоятельными оценками теоретической дисперсией, при этом D^* - смещенная оценка, а S^2 - несмещенная оценка

Заметим, что $D^* = \overline{X^2} - \overline{X}^2$

$$ED^* = E(\overline{X^2} - \overline{X}^2) = E\overline{X^2} - E(\overline{X}^2) = EX^2 - E(\overline{X}^2)$$

$$\text{Так как } D\overline{X} = E(\overline{X^2}) - (E\overline{X})^2, \text{ то } EX^2 - E(\overline{X}^2) = EX^2 - ((E\overline{X})^2 + D\overline{X}) = (EX^2 - EX) - D\overline{X} = \\ DX - D\overline{X} = DX - D\left(\frac{X_1 + \dots + X_n}{n}\right) = DX - \frac{1}{n^2} \sum_{i=1}^n DX_i = DX - \frac{1}{n^2} n DX_1 = DX - \frac{1}{n} DX = \frac{n-1}{n} DX,$$

то есть D^* - смещенная вниз оценка

$$ES^2 = E\left(\frac{n}{n-1} D^*\right) = \frac{n}{n-1} \frac{n-1}{n} DX = DX \implies S^2 - \text{несмещенная вниз оценка}$$

$$2. D^* = \overline{X^2} - \overline{X}^2 \xrightarrow{p} EX^2 - (EX)^2 = DX - \text{состоятельная оценка}$$

$$S^2 = \frac{n}{n-1} D^* \xrightarrow{p} DX$$

Nota. Отсюда видим, что выборочная дисперсия - асимптотически несмещенная оценка. Поэтому при большом (обычно не меньше 100) объеме выборке можно считать обычную выборочную дисперсию

Метод моментов (Пирсона)

Постановка задачи: пусть имеется выборка объема n неизвестного распределения, но известного типа, которое задается k параметрами: $\theta = (\theta_1, \theta_2, \dots, \theta_k)$. Требуется дать оценки данным неизвестным параметрам

Идея метода состоит в том, что сначала находим оценки k моментов, а затем с помощью теоретических формул из теории вероятности даем оценки этих параметров

Пусть \vec{X} - выборка из абсолютно непрерывного распределения F_θ с плотностью известного типа, которая задается k параметрами $f_\theta(x, \theta_1, \dots, \theta_k)$

Тогда теоретические моменты находим по формуле $m_i = \int_{-\infty}^{\infty} x^i f_\theta(x, \theta_1, \dots, \theta_k) dx = h_i(\theta_1, \dots, \theta_k)$

Получаем систему из k уравнений с k неизвестными. В эти уравнения подставляем найденные оценки моментов и, решая получившуюся систему уравнений, находим нужные оценки параметров

$$\begin{cases} \bar{x} = h_1(\theta_1^*, \dots, \theta_k^*) \\ \overline{x^2} = h_2(\theta_1^*, \dots, \theta_k^*) \\ \dots \\ \overline{x^k} = h_k(\theta_1^*, \dots, \theta_k^*) \end{cases}$$

Nota. Оценки по методу моментов как правило состоятельные, но часто смещенные

Ex. Пусть $X \in U(a, b)$. Обработав статданные, нашли оценки первого и второго моментов:

$$\bar{x} = 2.25; \overline{x^2} = 6.75$$

Найти оценки параметров a^*, b^*

$$\text{Плотность равномерного распределения } f_{(a,b)}(x) = \begin{cases} 0, & x < a \\ \frac{1}{b-a} & a \leq x \leq b, \\ 0, & x > b \end{cases}$$

$$EX = \int_a^b x \frac{1}{b-a} dx = \frac{a+b}{2}$$

$$EX = \int_a^b x^2 \frac{1}{b-a} dx = \frac{a^2 + ab + b^2}{3}$$

Получаем:

$$\begin{cases} \bar{x} = \frac{a^*+b^*}{2} \\ \bar{x^2} = \frac{a^{*2}+a^*b^*+b^{*2}}{3} \end{cases} \iff \begin{cases} \frac{a^*+b^*}{2} = 4.5 \\ a^{*2} + a^*b^* + b^{*2} = 20.25 \end{cases} \iff \begin{cases} \frac{a^*+b^*}{2} = 4.5 \\ a^*b^* = 0 \end{cases} \iff \begin{cases} a^* = 0 \\ b^* = 4.5 \end{cases}$$

Лекция 3.

Метод максимального правдоподобия

Пусть имеется выборка $\vec{X} = (X_1, \dots, X_n)$ из распределения известного типа, определяемого неизвестными параметрами $\theta = (\theta_1, \dots, \theta_n)$

Идея метода состоит в следующем: подбираем параметры таким образом, чтобы вероятность получения данной выборки при случайном эксперименте была наибольшей.

Если распределение дискретное, то $P_\theta(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P(X_1 = x_1) \dots P(X_n = x_n)$

Def. Функцией правдоподобия $L(\vec{X}, \theta)$ называется функция $L(\vec{X}, \theta) = P(X_1 = x_1) \dots P(X_n = x_n) = \prod_{i=1}^n P(X_i = x_i)$ при дискретном распределении

и $L(\vec{X}, \theta) = f_\theta(x_1) \dots f_\theta(x_n) = \prod_{i=1}^n f_\theta(x_i)$ в абсолютно непрерывном распределении

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

Nota. Так как $y = \ln x$ возрастающая функция, точки максимума совпадают, а такую функцию правдоподобия становится легче дифференцировать

Def. Оценкой максимального правдоподобия $\hat{\theta}$ называется значение θ , при котором функция правдоподобия $L(\vec{X}, \theta)$ достигает наибольшего значения (при фиксированных значениях выборки)

Ex. 1. Пусть $\vec{X} = (X_1, \dots, X_n)$ - выборка из распределения Пуассона Π_λ с неизвестным $\lambda > 0$

Мет. Для распределения Пуассона $P(X = x_i) = \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}$

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

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

$$\frac{\partial \ln L}{\partial \lambda} = \frac{n\bar{x}}{\lambda} - n = 0 \implies \hat{\lambda} = \bar{x} - \text{оценка максимального правдоподобия}$$

Убедимся, что этот экстремум - максимум: $\frac{\partial^2 \ln L}{\partial \lambda^2} = -\frac{n\bar{x}}{\lambda} < 0 \implies \hat{\lambda} = \bar{x} - \text{точка максимума}$

Ех. 2. Пусть (X_1, \dots, X_n) из $N(a, \sigma^2)$

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

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

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

$$\frac{\partial \ln L}{\partial a} = -\frac{1}{2\sigma^2} \sum_{i=1}^n -2(x_i - a) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - a) = \frac{n\bar{x} - na}{\sigma^2}$$

$$\frac{\partial \ln L}{\partial \sigma} = -\frac{n}{\sigma} - \sum_{i=1}^n (x_i - a)^2 \frac{1}{2} \cdot (-2) \cdot \sigma^{-3} = \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - a)^2 - \frac{n}{\sigma}$$

$$\begin{cases} \frac{n\bar{x} - na}{\sigma^2} = 0 \\ \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - a)^2 - \frac{n}{\sigma} = 0 \end{cases} \implies \begin{cases} \hat{a} = \bar{x} \\ \widehat{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (x_i - a)^2 = D^* \end{cases}$$

Ех. 3. Пусть (X_1, \dots, X_n) из $U(0, \theta)$. Найти оценку θ этого распределения.

Воспользуемся методом моментов:

$$EX = \frac{a+b}{2} = \frac{\theta}{2} \implies \bar{x} = \frac{\theta^*}{2} \implies \theta^* = 2\bar{x}$$

Воспользуемся методом максимального правдоподобия:

$$f_\theta = \begin{cases} 0, & x < 0 \\ \frac{1}{\theta}, & 0 \leq x \leq \theta \\ 0, & x > \theta \end{cases}$$

$$X_{(n)} = \max_i (X_1, \dots, X_n)$$

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

$L(\vec{X}, \theta)$ достигает наибольшего значения при наименьшем значении θ^n , то есть при $\hat{\theta} = X_{(n)}$

Сравним оценки:

$$\theta^* = 2\bar{x} - \text{несмещенная оценка, так как } E\theta^* = 2E\bar{x} = 2EX = \theta$$

$$E(\theta^* - \theta)^2 = D\theta^* = D2\bar{x} = 4D\bar{x} = 4 \frac{D\bar{x}}{n} = \frac{4}{n} \frac{\theta^2}{12} = \frac{\theta^2}{3n}$$

Изучим распределение $X_{(n)}$: $F_{X_{(n)}}(x) = P(X_{(n)} < x) = P(X_1 < x, X_2 < x, \dots, X_n < x) = P(X_1 < x) \dots P(X_n < x) = F_{X_1}(x) \dots F_{X_n}(x) = F_{(x_1)}^n(x)$

$$F_{X_1}(x) = \begin{cases} 0, & x < 0 \\ \frac{x}{\theta}, & 0 \leq x \leq \theta \\ 1, & x > \theta \end{cases} \Rightarrow F_{X_{(n)}}(x) = \begin{cases} 0, & x < 0 \\ \frac{x^n}{\theta^n}, & 0 \leq x \leq \theta \\ 1, & x > \theta \end{cases} \Rightarrow f_{X_{(n)}}(x) = \begin{cases} 0, & x < 0 \\ n \frac{x^{n-1}}{\theta^n}, & 0 \leq x \leq \theta \\ 1, & x > \theta \end{cases}$$

$$EX_{(n)} = \int_0^\theta x \cdot \frac{nx^{n-1}}{\theta^n} dx = \frac{n}{\theta^n} \int_0^\theta x^n dx = \frac{nx^{n+1}}{\theta^n(n+1)} \Big|_0^\theta = \frac{n\theta}{n+1} - \text{смещенная вниз оценка}$$

$\tilde{\theta} = \frac{n+1}{n} X_{(n)}$ - несмещенная оценка (будем считать, что эффективность не изменилась)

$$E\tilde{\theta}^2 = E\left(\frac{n+1}{n} X_{(n)}\right)^2 = \frac{(n+1)^2}{n^2} EX_{(n)}^2 = \frac{(n+1)^2}{n^2} \int_0^\theta x^2 \frac{nx^{n-1}}{\theta^n} dx = \frac{(n+1)^2 \theta^2}{n(n+2)}$$

$$D\tilde{\theta} = E\tilde{\theta}^2 - (E\tilde{\theta})^2 = \frac{\theta^2}{n(n+2)}$$

$$D\tilde{\theta} = \frac{\theta^2}{n(n+2)} < \frac{\theta^2}{3n} = D\theta^*$$

Таким образом, оценка по методу правдоподобия сходится быстрее, чем оценка по методу моментов, поэтому она лучше

Отсюда следует, что при равномерном распределении выборочное среднее не является эффективной оценкой для математического ожидания; вместо нее половина максимального элемента выборки будет лучше

Nota. Эффективной здесь будет несмещенная оценка $\frac{n+1}{2n} X_{(n)}$

В общем случае для $U(a, b)$ будет такая эффективная оценка математического ожидания - $\frac{X_{(1)} + X_{(n)}}{2}$, длины интервала - $\frac{n+1}{n-1} (X_{(n)} - X_{(1)})$

Nota. При методе максимального правдоподобия обычно получаем состоятельные и эффективные оценки, но часто смещенные

Неравенство Рао-Крамера

Пусть $X \in F_\theta$ - семейство распределений с параметром $\theta \in \mathbb{R}$

Def. Носителем семейства распределений F_θ называется множество $C \subset \mathbb{R}$ такое, что $P(X \in C) = 1 \ \forall X \in F_\theta$

$$f_\theta(x) = \begin{cases} \text{плотность } f_\theta(x) \text{ при непрерывном распределении} \\ P_\theta(X = x) \text{ при дискретном распределении} \end{cases}$$

Def. Информацией Фишера $I(\theta)$ семейства распределений F_θ называется величина $I(\theta) =$

$E \left(\frac{\partial}{\partial \theta} \ln f_{\theta}(X) \right)^2$ при условии, что она существует

Def. Семейство распределений F_{θ} называется регулярным, если:

- существует носитель C семейства F_{θ} такой, что $\forall x \in C$ функция $\ln f_{\theta}(x)$ непрерывно дифференцируема по θ
- информация Фишера $I(\theta)$ существует и непрерывна по θ

Th. Пусть (X_1, \dots, X_n) - выборка объема n из регулярного семейства F_{θ} ,
 $\theta^* = \theta^*(X_1, \dots, X_n)$ - несмещенная оценка параметра θ , дисперсия которой $D\theta^*$ ограничена
 в любой замкнутой ограниченной области параметра θ

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

Следствие: если при данных условиях получили $D\theta^* = \frac{1}{nI(\theta)}$, то оценка θ^* является эффективной (то есть дальше улучшать уже некуда)

Ex. Пусть (X_1, \dots, X_n) из $N(a, \sigma^2)$ (то есть $F_a = N(a, \sigma^2)$, σ^2 зафиксируем)

Проверим эффективность $a^* = \bar{x}$

Плотность $f_a(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-a)^2}{2\sigma^2}}$, носитель - вся прямая \mathbb{R}

$$\ln f_a(x) = -\ln \sigma - \frac{1}{2} \ln 2\pi - \frac{(x-a)^2}{2\sigma^2}, \quad a \in (-\infty, \infty)$$

$$\frac{\partial}{\partial a} \ln f_a(x) = \frac{1}{2\sigma^2} \cdot 2(x-a) = \frac{x-a}{\sigma^2} - \text{непрерывна для всех } a \in \mathbb{R}$$

$$I(a) = E \left(\frac{\partial}{\partial a} \ln f_a(X) \right)^2 = E \left(\frac{X-a}{\sigma^2} \right)^2 = \frac{1}{\sigma^4} E(X-a)^2 = \frac{E(X-EX)^2}{\sigma^4} = \frac{DX}{\sigma^4} = \frac{1}{\sigma^2} - \text{непрерывна по } a$$

Из этого следует, что $N(a, \sigma^2)$ - регулярное семейство относительно параметра a

$$Da^* = D\bar{x} = \frac{DX}{n} = \frac{\sigma^2}{n} - \text{ограничена по параметру } a$$

По неравенству Рао-Крамера $Da^* = \frac{\sigma^2}{n} = \frac{1}{nI(a)} = \frac{1}{n} \sigma^2$; из этого следует, что a^* - эффективная оценка параметра a

Nota. Аналогично можно показать, что S^2 - несмещенная эффективная оценка для параметра σ^2

Лекция 4.

Основные распределения математической статистики

Def. Случайная величина имеет нормальное распределение $\xi \in N(a, \sigma^2)$ с параметрами a и σ^2 , если ее плотность имеет вид $f_\xi(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-a)^2}{2\sigma^2}}$

На практике нормальное распределение встречается чаще всего в силу ЦПТ

Def. Распределение $N(0, 1)$ с параметрами $a = 0, \sigma^2 = 1$ называется стандартным нормальным распределением. Его плотность равна $\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$. В дальнейшем такую случайную величину будем называть стандартной нормалью

Свойства

1. $a = E\xi \quad \sigma^2 = D\xi$
2. Линейность: $\xi \in N(a, \sigma^2)$, то $\eta = b\xi + \gamma \in N(ab + \gamma, b^2\sigma^2)$
3. Стандартизация: Если $\xi \in N(a, \sigma^2)$, то $\eta = \frac{\xi - a}{\sigma} \in N(0, 1)$
4. Устойчивость относительно суммирования: если $\xi_1 \in N(a_1, \sigma_1^2)$, $\xi_2 \in N(a_2, \sigma_2^2)$, независимы то $\xi_1 + \xi_2 \in N(a_1 + a_2, \sigma_1^2 + \sigma_2^2)$

Распределение «хи-квадрат»

Def. Распределение «хи-квадрат» H_n со степенями свободы n называется распределение суммы квадратов независимых стандартных нормальных величин: $\chi_n^2 = X_1^2 + X_2^2 + \dots + X_n^2$, где $X \in N(0, 1)$ и независимы

Свойства

1. $E\chi_n^2 = n$

Так как $\forall i \ X_i \in N(0, 1)$, то $E\chi_n^2 = D\chi_n^2 + (E\chi_n)^2 = 1 \implies E(X_1^2 + \dots + X_n^2) = \sum_{i=1}^n E X_i^2 = n$

2. Устойчивость относительно суммирования: если $X \in H_n$, $Y \in H_m$, независимы, то $X + Y \in H_{n+m}$ (по определению)
3. $\frac{\chi_k^2}{k} \xrightarrow[k \rightarrow \infty]{p} 1$ (по Закону Больших Чисел)

Распределение Стюдента

Def. Пусть X_0, X_1, \dots, X_k - независимые стандартные нормальные величины. Распределением Стюдента T_k с k степенями свободы называется распределение случайной величины $t_k =$

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

Свойства

1. $Et_k = 0$ - в силу симметрии
2. $t_k \Rightarrow N(0, 1)$ (на практике при $k \geq 100$ распределение Стюдента можно считать стандартным нормальным)

Распределение Фишера-Снедекера

Def. Распределением Фишера-Снедекера $F_{n,m}$ (другое название - F-распределение) со степенями свободы n и m называется распределение случайной величины $f_{n,m} = \frac{\frac{\chi_n^2}{n}}{\frac{\chi_m^2}{m}}$, где χ_n^2 и χ_m^2 - независимые случайные величины с распределением «хи-квадрат»

Свойства

1. $Ef_{n,m} = \frac{n}{n-2}$
2. $f_{n,m} \xrightarrow[n, m \rightarrow \infty]{p} 1$

Математическое ожидание и дисперсия случайного вектора

Пусть $\vec{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$ - случайный вектор, где случайная величина X_i - компонента (координата) случайного вектора

Def. Математическим ожиданием случайного вектора называется вектор с координатами из математических ожиданий компонент: $E\vec{X} = \begin{pmatrix} EX_1 \\ \vdots \\ EX_n \end{pmatrix}$

Def. Дисперсией случайного вектора (или матрицей ковариаций) случайного вектора \vec{X} называется матрица $D\vec{X} = E(\vec{X} - E\vec{X})(\vec{X} - E\vec{X})^T$, состоящая из элементов $d_{ij} = \text{cov}(X_i, X_j)$

Nota. На главной диагонали стоят дисперсии компонент: $d_{ii} = DX_i$

Nota. $D\vec{X}$ - симметричная положительно определенная матрица

Свойства

1. $E(A\vec{X}) = AE\vec{X}$
2. $E(\vec{X} + \vec{B}) = E\vec{X} + \vec{B}$, где \vec{B} - вектор чисел
3. $D(A\vec{X}) = A \cdot D\vec{X} \cdot A^T$
4. $D(\vec{X} + \vec{B}) = D\vec{X}$

Многомерное нормальное распределение

Def. Пусть случайный вектор $\vec{\xi} = \begin{pmatrix} \xi_1 \\ \vdots \\ \xi_n \end{pmatrix}$ имеет вектор средних $\vec{a} = E\vec{\xi}$, K - симметричная положительно определенная матрица. Вектор $\vec{\xi}$ имеет нормальное распределение в \mathbb{R}^n с параметрами \vec{a} и K , если его плотность $f_{\vec{\xi}}(\vec{X}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{\det K}} e^{-\frac{1}{2}(\vec{X}-\vec{a})^T K^{-1}(\vec{X}-\vec{a})}$

Свойства

1. Матрица $K = D\vec{\xi} = (\text{cov}(\xi_i, \xi_j))$ - матрица ковариаций
2. При $\vec{a} = \vec{0}$ и $K = E$ имеем вектор из независимых стандартных нормальных величин

$$\text{При } \vec{a} = \vec{0} \text{ и } K = E: f_{\vec{\xi}}(X_1, \dots, X_n) = \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2} \begin{pmatrix} X_1 & \dots & X_n \end{pmatrix} E \begin{pmatrix} X_1 & \dots & X_n \end{pmatrix}^T} =$$

$$\frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2}(X_1^2 + \dots + X_n^2)} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}X_1^2} + \dots + \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}X_n^2}$$

Так как плотность распалась на произведение плотностей стандартного нормального распределения, то все компоненты имеют стандартное нормальное распределение

Далее вектором из независимых стандартных нормальных величин для краткости будем называть стандартным нормальным вектором

3. $\exists \vec{X}$ - стандартный нормальный вектор, B - невырожденная матрица, тогда вектор $\vec{Y} = B\vec{X} + \vec{a}$ имеет многомерное нормальное распределение с параметрами \vec{a} и $K = BB^T$
4. $\exists \vec{Y} \in N(\vec{a}, K)$. Тогда вектор $\vec{X} = B^{-1}(\vec{Y} - \vec{a})$ - стандартный нормальный вектор, где $B = \sqrt{K}$

Следствие. Эквивалентное определение: Многомерное нормальное распределение - это то, которое получается из стандартного нормального вектора при помощи невырожденного преобразования и сдвиг

5. $\exists \vec{X}$ - стандартный нормальный вектор, C - ортогональная матрица. Тогда $\vec{Y} = C\vec{X}$ - стандартный нормальный вектор

Так как C - ортогональная, то $C^T = C^{-1}$. Тогда по третьему свойству $K = CC^T = E$, а по второму свойству \vec{Y} - стандартный нормальный вектор

6. \square случайный вектор $\xi \in N(\vec{a}, K)$. Тогда его координаты независимы тогда и только тогда, когда они не коррелированы (то есть матрица ковариаций K диагональная)

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

Многомерная центральная предельная теорема

Th. Среднее арифметическое независимых одинаково распределенных случайных векторов слабо сходится к многомерному нормальному распределению

Лемма Фишера

Пусть вектор \vec{X} - стандартный нормальный вектор, C - ортогональная матрица, $\vec{Y} = C\vec{X}$. Тогда $\forall 1 \leq k \leq n-1$ случайная величина $T(\vec{X}) = \sum_{i=1}^n X_i^2 - Y_1^2 - Y_2^2 - \dots - Y_k^2$ не зависит от Y_1, Y_2, \dots, Y_k и имеет распределение «хи-квадрат» со степенями свободы $n-k$

Так как C - ортогональное преобразование, то $\|\vec{X}\| = \|\vec{Y}\|$, то есть $\sum_{i=1}^n X_i^2 = \sum_{i=1}^n Y_i^2 \implies$

$$T(\vec{X}) = \sum_{i=1}^n X_i^2 - Y_1^2 - Y_2^2 - \dots - Y_k^2 = Y_{k+1}^2 + \dots + Y_n^2$$

Согласно свойству 5 $Y_i \in N(0, 1)$ и независимы, то по определению «хи-квадрат» $T(\vec{X}) \in H_{n-k}$ и не зависит от Y_1, \dots, Y_k

Основная теорема

Эта теорема также известна как [основное следствие леммы Фишера](#)

Th. Пусть (X_1, \dots, X_n) - выборка из нормального распределения $N(a, \sigma^2)$, \bar{x} - выборочное среднее, S^2 - исправленная дисперсия.

Тогда справедливы следующие высказывания:

1. $\sqrt{n} \frac{\bar{x} - a}{\sigma} \in N(0, 1)$
2. $\sum_{i=1}^n \frac{(X_i - a)^2}{\sigma^2} \in H_n$
3. $\sum_{i=1}^n \frac{(X_i - \bar{x})^2}{\sigma^2} = \frac{nD^*}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2} \in H_{n-1}$
4. $\sqrt{n} \frac{\bar{x} - a}{S} \in T_{n-1}$
5. \bar{x} и S^2 независимы

1. Так как $X_i \in N(a, \sigma^2)$, то $\sum_{i=1}^n X_i \in N(na, n\sigma^2) \Rightarrow \bar{x} \in N\left(a, \frac{\sigma^2}{n}\right) \Rightarrow \bar{x} - a \in N\left(0, \frac{\sigma^2}{n}\right) \Rightarrow \frac{\sqrt{n}}{\sigma}(\bar{x} - a) \in N(0, 1)$

2. Так как $X_i \in N(a, \sigma^2)$, то $\frac{X_i - a}{\sigma} \in N(0, 1)$ и $\sum_{i=1}^n \frac{(X_i - a)^2}{\sigma^2} \in H_n$ по определению

3. $\sum_{i=1}^n \frac{(X_i - \bar{x})^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{X_i - a}{\sigma} - \frac{\bar{x} - a}{\sigma} \right)^2 = \sum_{i=1}^n (z_i - \bar{z})^2$, где $z_i = \frac{X_i - a}{\sigma} \in N(0, 1)$, $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i = \frac{\sum_{i=1}^n X_i - na}{\sigma} = \frac{\bar{x} - a}{\sigma}$

Поэтому можно считать, что изначально $X_i \in N(0, 1)$

$T(\vec{X}) = \sum_{i=1}^n (X_i - \bar{x})^2 = nD^* = n(\bar{x}^2 - \bar{x}^2) = \sum_{i=1}^n X_i^2 - n\bar{x}^2 = \sum_{i=1}^n X_i^2 - Y_1^2$, где $Y_1 = \sqrt{n}\bar{x} = \frac{X_1}{\sqrt{n}} + \dots + \frac{X_n}{\sqrt{n}}$

Строчка $\left(\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}}\right)$ имеет длину 1, поэтому ее можно дополнить до ортогональной матрицы C , тогда Y_1 - первая компонента $\vec{Y} = C\vec{X}$, и согласно лемме Фишера

$T(\vec{X}) = \sum_{i=1}^n (X_i - \bar{x})^2 \in H_{n-1}$

5. Согласно лемме Фишера $T(\vec{X}) = (n-1)S^2$ не зависит от $Y_1 = \sqrt{n}\bar{x} \Rightarrow S^2$ и \bar{x} - независимы

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

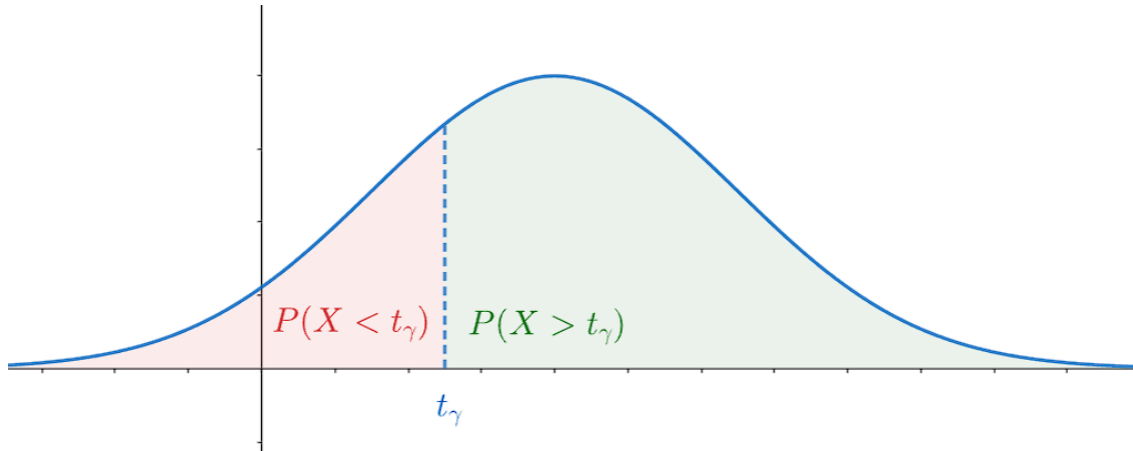
Так как по пятому пункту числитель и знаменатель независимы, по определению получаем распределение Стьюдента

Лекция 5.

Квантильное распределение

Предполагаем, что распределение абсолютно непрерывное и $F(x)$ - функция распределения

Def. 1. Число t_γ называется квантилем распределения уровня γ , если значения функции распределения $F(t_\gamma) = \gamma$ или $P(X < t_\gamma) = \gamma$ ($t_\gamma = F^{-1}(\gamma)$)



Ех. Медиана - квантиль уровня $\frac{1}{2}$

Def. 2. Число t_α называется квантилем уровня значимости α , если $P(X > t_\alpha) = \alpha$ или $F(t_\alpha) = 1 - \alpha$
 Ясно, что $\gamma = 1 - \alpha$

Интервальные оценки

Недостатки точечных оценок - неизвестно насколько они далеки от реального значения параметра и насколько им можно доверять. Особенно это заметно при малых выборках. Поэтому мы указываем интервал, в котором лежит этот параметр с заданной вероятностью (надежностью) γ . Такие оценки называются интервальными (доверительными)

Def. Интервал $(\theta_\gamma^-; \theta_\gamma^+)$ называется доверительным интервалом параметра θ надежности γ , если вероятность $P(\theta_\gamma^- < \theta < \theta_\gamma^+) = \gamma$

Nota. Если имеем дискретную случайную величину, то $P(\theta_\gamma^- < \theta < \theta_\gamma^+) \geq \gamma$

Nota. Так как параметр θ - константа, то бессмысленно говорить о его попадании в интервал.

Правильно: доверительный интервал накрывает параметр θ с вероятностью γ

Nota. 1. $\alpha = 1 - \gamma$ называется уровнем значимости доверительного интервала

Nota. 2. Обычно пытаются строить симметричный доверительный интервал относительно несмещенной оценки θ^*

Nota. 3. Возникает вопрос, какой уровень γ выбрать для исследования. Стандартные уровни надежности γ : 0.9, 0.95, 0.99, 0.999. Самый мейнстримный - 0.95. В малых выборках используют 0.9

Вспомним основную теорему:

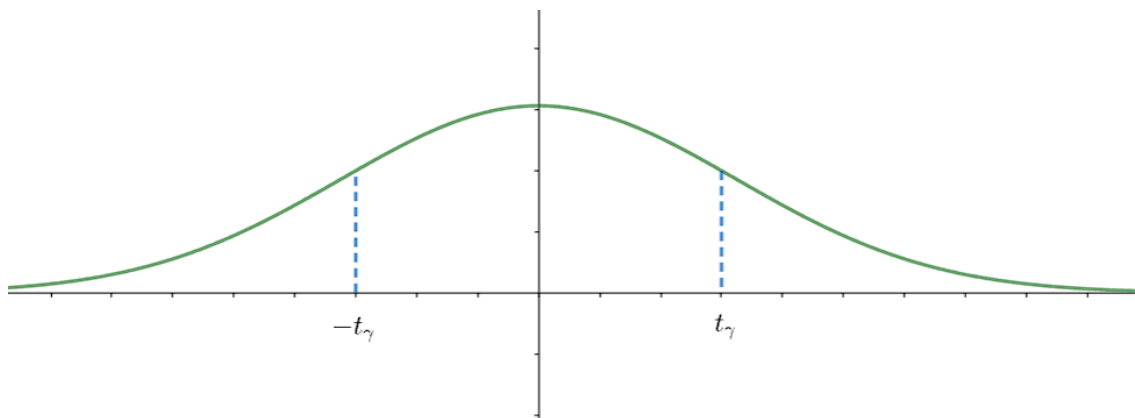
$\square (X_1, \dots, X_n)$ - выборка объема n из $N(\alpha, \sigma^2)$

\bar{x} - выборочное среднее, S^2 - исправленная дисперсия, D^* - выборочная дисперсия

Тогда:

1. $\sqrt{n} \frac{\bar{x} - a}{\sigma} \in N(0, 1)$
2. $\sum_{i=1}^n \left(\frac{X_i - a}{\sigma} \right)^2 = \frac{n\tilde{\sigma}^2}{\sigma^2} \in H_n$, где $n\tilde{\sigma}^2 = \sum_{i=1}^n (X_i - a)^2$
3. $\sum_{i=1}^n \left(\frac{X_i - \bar{x}}{\sigma} \right)^2 = \frac{(n-1)S^2}{\sigma^2} = \frac{nD^*}{\sigma^2} \in H_{n-1}$
4. $\sqrt{n} \frac{\bar{x} - a}{S} \in T_{n-1}$
5. \bar{x} и S^2 - независимы

Nota. Если $F(x)$ - функция симметричного относительно $x = 0$ распределения, то $P(|X| < t) = 2F(t) - 1$



Доверительные интервалы для параметров нормального распределения

Пусть $\vec{X} = (X_1, \dots, X_n)$ - выборка объема n из $N(a, \sigma^2)$. Хотим найти интервалы для параметров a и σ^2

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

По пункту 1 из теоремы $\sqrt{n} \frac{\bar{x} - a}{\sigma} \in N(0, 1)$

$$P\left(-t_\gamma < \sqrt{n} \frac{\bar{x} - a}{\sigma} < t_\gamma\right) = P\left(\left|\sqrt{n} \frac{\bar{x} - a}{\sigma}\right| < t_\gamma\right) = 2F_0(t_\gamma) - 1 = \gamma$$

$$F_0(t_\gamma) = \frac{1+\gamma}{2} \implies t_\gamma - \text{квантиль уровня } \frac{1+\gamma}{2} \text{ для } N(0, 1), \text{ где } F_0(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{z^2}{2}} dz$$

Решая неравенство, получаем $-t_\gamma < \sqrt{n} \frac{\bar{x} - a}{\sigma} < t_\gamma$

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

$$\bar{x} - t_\gamma \frac{\sigma}{\sqrt{n}} < a < \bar{x} + t_\gamma \frac{\sigma}{\sqrt{n}} - \text{симметричный интервал относительно } \bar{x}$$

Доверительный интервал надежности γ : $\left(\bar{x} - t_\gamma \frac{\sigma}{\sqrt{n}}, \bar{x} + t_\gamma \frac{\sigma}{\sqrt{n}}\right)$, где t_γ - квантиль $N(0, 1)$

уровня $\frac{1+\gamma}{2}$

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

Из пункта 4 из теоремы $\sqrt{n} \frac{\bar{x} - a}{S} \in T_{n-1}$

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

$$F_{T_{n-1}}(t_\gamma) = \frac{1+\gamma}{2} \implies t_\gamma - \text{квантиль } T_{n-1} \text{ уровня } \frac{1+\gamma}{2}$$

Аналогично с примером выше получаем интервал $\left(\bar{x} - t_\gamma \frac{S}{\sqrt{n}}, \bar{x} + t_\gamma \frac{S}{\sqrt{n}}\right)$, где t_γ - квантиль

T_{n-1} уровня $\frac{1+\gamma}{2}$

III. Доверительный интервал для параметра σ^2 при неизвестном a

По пункту 3 из теоремы $\sum_{i=1}^n \left(\frac{X_i - \bar{x}}{\sigma}\right)^2 = \frac{(n-1)S^2}{\sigma^2} = \frac{nD^*}{\sigma^2} \in H_{n-1}$

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

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

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

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

$$\frac{(n-1)S^2}{\chi_2^2} < \sigma^2 < \frac{(n-1)S^2}{\chi_1^2} \text{ или } \frac{nD^*}{\chi_2^2} < \sigma^2 < \frac{nD^*}{\chi_1^2}$$

Получаем 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 - квантили H_{n-1} уровней $\frac{1-\gamma}{2}$ и

$\frac{1+\gamma}{2}$

Nota. Данный интервал не симметричен относительно неизвестного параметра σ^2

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

По пункту 2 из теоремы $\sum_{i=1}^n \left(\frac{X_i - a}{\sigma} \right)^2 = \frac{n\tilde{\sigma}^2}{\sigma^2} \in H_{n-1}$

Пусть χ_1^2 и χ_2^2 - квантили H_n уровней $\frac{1-\gamma}{2}$ и $\frac{1+\gamma}{2}$

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

Аналогично получаем интервал $\left(\frac{n\tilde{\sigma}^2}{\chi_2^2}, \frac{n\tilde{\sigma}^2}{\chi_1^2}\right)$, где χ_1^2 и χ_2^2 - квантили H_n уровней $\frac{1-\gamma}{2}$ и

$$\frac{1+\gamma}{2}, n\tilde{\sigma}^2 = \sum_{i=1}^n (X_i - a)^2$$

$$\text{Nota. } \tilde{\sigma}^2 - D^* = \frac{1}{n} \sum_{i=1}^n (X_i - a)^2 - \frac{1}{n} \sum_{i=1}^n (X_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n (X_i^2 - 2aX_i + a^2 - X_i^2 + 2\bar{x}X_i - \bar{x}^2) = \frac{1}{n} (na^2 -$$

$$2an\bar{x} + 2\bar{x} \cdot n\bar{x} - n\bar{x}^2) = a^2 - 2a\bar{x} + \bar{x}^2 = (a - \bar{x})^2 \implies \tilde{\sigma}^2 = D^* + (a - \bar{x})^2$$

$$\text{Получаем } \left(\frac{n(D^* + (a - \bar{x})^2)}{\chi_2^2}, \frac{n(D^* + (a - \bar{x})^2)}{\chi_1^2} \right)$$

Асимптотические доверительные интервалы

Def. Интервал $(\theta_\gamma^-, \theta_\gamma^+)$ называется асимптотическим доверительным интервалом параметра θ уровня γ , если $P(\theta_\gamma^- < \theta < \theta_\gamma^+) \xrightarrow{n \rightarrow \infty} \gamma$

Ex. Доверительный интервал вероятности события A

Пусть $p = P(A)$, $q = 1 - p$, n - число испытаний или объем выборки (X_1, \dots, X_n) , где $X_i \in \{0, 1\}$

$$p^* = \frac{n_A}{n} = \bar{x} - \text{оценка } p$$

Согласно Центральной предельной теореме $\sqrt{n} \frac{p^* - p}{DX_1} = \sqrt{n} \frac{p^* - p}{\sqrt{pq}} \Rightarrow N(0, 1)$

Так как $p^* \xrightarrow{p} p$, то $\sqrt{n} \frac{p^* - p}{\sqrt{p^*(1-p^*)}} = \underbrace{\sqrt{n} \frac{p^* - p}{\sqrt{p(1-p)}}}_{\Rightarrow N(0,1)} \underbrace{\frac{\sqrt{p(1-p)}}{\sqrt{p^*(1-p^*)}}}_{\xrightarrow{p} 1} \Rightarrow N(0, 1)$

$$P\left(\left|\sqrt{n} \frac{p^* - p}{\sqrt{p^*(1-p^*)}}\right| < t_\gamma\right) \xrightarrow{n \rightarrow \infty} 2F_0(t_\gamma) - 1 = \gamma$$

$$F_0(t_\gamma) = \frac{1+\gamma}{2}, t_\gamma - \text{квантиль } N(0, 1) \text{ уровня } \frac{1+\gamma}{2}$$

$$\text{Получаем } \left|\sqrt{n} \frac{p^* - p}{\sqrt{p^*(1-p^*)}}\right| < t_\gamma$$

$$|p^* - p| < t_\gamma \frac{\sqrt{p^*(1-p^*)}}{\sqrt{n}}$$

Итак, $\left(-t_\gamma \frac{\sqrt{\bar{x}(1-\bar{x})}}{\sqrt{n}}, t_\gamma \frac{\sqrt{\bar{x}(1-\bar{x})}}{\sqrt{n}}\right)$, где t_γ - квантиль $N(0, 1)$ уровня $\frac{1+\gamma}{2}$

Лекция 6.

Проверка статистических гипотез

$\vec{X} = (X_1, \dots, X_n)$ из некоторого распределения F

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

Доказать какое-то утверждение с помощью методов матстатистики невозможно - можно лишь с какой-то долей уверенности утверждать

Def. Гипотеза называется простой, если она однозначно определяет распределение: $H : F = F_1$, где F_1 - распределение известного типа с известными параметрами

В противном случае гипотеза называется сложной - она является объединением конечного или бесконечного числа гипотез

Например, «величина X принадлежит нормальному распределению» - сложная гипотеза, а «величина X принадлежит нормальному распределению с матожиданием $a = 1$ и дисперсией $\sigma^2 = 1$ » - простая

В общем случае работаем со схемой из двух или более гипотез. В ходе проверки принимается ровно одна из них. Мы ограничимся самой простой схемой из 2 гипотез: H_0 - основная (нулевая) гипотеза, $H_1 = \bar{H}_0$ - альтернативная (конкурирующая) гипотеза, состоящая в том, что основная гипотеза неверна

Основная гипотеза H_0 принимается или отклоняется при помощи статистики критерия K

$K(X_1, \dots, X_n) \longrightarrow \mathbb{R} = \bar{S} \cup S \longrightarrow (H_0, H_1)$

$$\begin{cases} H_0, & \text{если } K(X_1, \dots, X_n) \in \bar{S} \\ H_1, & \text{если } K(X_1, \dots, X_n) \in S \end{cases}$$

Вместо «гипотеза доказана» лучше употреблять «гипотеза принимается/отвергается»

Область S называется критической областью, а точка $t_{кр}$ на границе областей называется критической

Def. Ошибка первого рода состоит в том, что H_0 отклоняется, хотя она верна. Аналогично, ошибка второго рода состоит в том, что H_1 отклоняется, хотя она верна.

Def. Вероятность α ошибки первого рода называется уровнем значимости критерия. Вероятность ошибки второго рода обозначаем β . Мощностью критерия называется вероятность $1 - \beta$ (вероятность недопущения ошибки второго рода)

Ясно, что критерий будет тем лучше, чем меньше вероятности ошибок α и β . При увеличении объема выборки уменьшаются обе вероятности. При фиксированном объеме попытки уменьшить одну вероятность увеличат другую

Одним из способов является фиксация одной вероятности (принято α) и уменьшение другой

Построение критериев согласия

Def. Говорят, что критерий K является критерием асимптотического уровня ε , если вероятность ошибки первого рода $\alpha \xrightarrow[n \rightarrow \infty]{} \varepsilon$

Def. Критерий K для проверки гипотезы H_0 называется состоятельным, если вероятность ошибки второго рода $\beta \xrightarrow[n \rightarrow \infty]{} 0$

Def. Критерием согласия уровня ε называем состоятельный критерий асимптотического уровня ε

Обычно критерий согласия строится по следующей схеме: берется статистика $K(X_1, \dots, X_n)$, обладающая свойствами:

1. Если H_0 верна, то $K(X_1, \dots, X_n) \Rightarrow Z$, где Z - известное распределение
2. Если H_0 неверна, то есть верна H_1 , то $K(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{p} \infty$ (достаточно сильно отклоняться от распределения Z)

Построенный таким образом критерий является критерием согласия, то есть обладает свойствами

1. критерия асимптотического уровня
2. состоятельного критерия

Пусть $t_{кр}$ - критическая точка такая, что $P(|Z| > t_{кр}) = \varepsilon$ - заданный уровень ошибки первого рода

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

1. Тогда $\alpha = P(|K| \geq t_{кр} \mid H_0) = 1 - P(|K| < t_{кр} \mid H_0) = 1 - (F_K(t_{кр}) - F_K(-t_{кр}))$
т.к. при верной H_0 $F_K(x) \xrightarrow[n \rightarrow \infty]{} F_Z(x)$

$$F_K(-t_{кр})) \xrightarrow[n \rightarrow \infty]{} 1 - (F_Z(t_{кр}) - F_Z(-t_{кр})) = P(|Z| \geq t_{кр}) = \varepsilon$$

2. Если H_1 верна, то $|K| \xrightarrow{p} \infty$, то есть $\forall C \ P(|K| > C \mid H_1) \xrightarrow{p} 1 \implies \beta = P(|K| < C \mid H_1) \xrightarrow{p} 0$

Гипотеза о среднем нормальной совокупности при известной дисперсии

$\vec{X} = (X_1, \dots, X_n)$ из $N(a, \sigma^2)$, причем σ^2 известен.

Проверяется гипотеза, что $H_0 : a = a_0$, против $H_1 : a \neq a_0$ для уровня значимости α

1. По пункту 1 теоремы, если $H_0 : a = a_0$ верна, то $K = \sqrt{n} \frac{\bar{x} - a_0}{\sigma} = \sqrt{n} \frac{\bar{x} - a}{\sigma} \in N(0, 1)$

2. Если верна $H_1 : a \neq a_0$, то $|K| = \sqrt{n} \left| \frac{\bar{x} - a_0}{\sigma} \right| = \sqrt{n} \left| \frac{\bar{x} - a}{\sigma} + \frac{a - a_0}{\sigma} \right| =$

$$= \left| \underbrace{\sqrt{n} \frac{\bar{x} - a}{\sigma}}_{\substack{\in N(0,1), \text{ограничен} \\ \text{по вероятности}}} + \underbrace{\sqrt{n} \frac{a - a_0}{\sigma}}_{\substack{\rightarrow \infty \\ \text{const}}} \right| \xrightarrow{p} \infty$$

Для уровня значимости α находим $t_{кр}$ такую, что $\alpha = P(|K| \geq t_{кр} \mid H_0) = P(|Z| \geq t_{кр}) \implies P(|Z| < t_{кр}) = 2F_0(t_{кр}) - 1 = 1 - \alpha$

$F_0(t_{кр}) = 1 - \frac{\alpha}{2}$ - то есть $t_{кр}$ - квантиль стандартного нормального распределения уровня $1 - \frac{\alpha}{2}$

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

Гипотеза о среднем нормальной совокупности при неизвестной дисперсии

1. По пункту 4 основной теоремы, если $H_0 : a = a_0$ верна, то $K = \sqrt{n} \frac{\bar{x} - a_0}{S} = \sqrt{n} \frac{\bar{x} - a}{S} \in T_{n-1}$

2. Если верна $H_1 : a \neq a_0$, то $|K| = \sqrt{n} \left| \frac{\bar{x} - a_0}{S} \right| = \sqrt{n} \left| \frac{\bar{x} - a}{S} + \frac{a - a_0}{S} \right| =$

$$= \left| \underbrace{\sqrt{n} \frac{\bar{x} - a}{S}}_{\substack{\in T_{n-1}, \text{ограничен} \\ \text{по вероятности}}} + \underbrace{\sqrt{n} \frac{a - a_0}{S}}_{\substack{\rightarrow \infty \\ \text{const}}} \right| \xrightarrow{p} \infty$$

Аналогично получаем $t_{кр}$ - квантиль распределения T_{n-1} уровня $1 - \frac{\alpha}{2}$

Доверительные интервалы как критерии гипотез по параметрам распределения

$\square(X_1, \dots, X_n)$ из F_θ , где F_θ - распределение известного типа с неизвестным параметром θ

Проверяется гипотеза, что $H_0 : \theta = \theta_0$, против $H_1 : \theta \neq \theta_0$

Допустим, что для θ построен доверительный интервал (θ_Y^-, θ_Y^+) , то есть $P(\theta_Y^- < \theta < \theta_Y^+) = \gamma$.

Тогда критерий $\begin{cases} H_0, & \text{если } \theta_0 \in (\theta_Y^-, \theta_Y^+) \\ H_1, & \text{если } \theta_0 \notin (\theta_Y^-, \theta_Y^+) \end{cases}$ будет уровня $\alpha = 1 - \gamma$

$$\alpha = P(\theta_0 \notin (\theta_Y^-, \theta_Y^+) \mid H_0) = 1 - P(\theta_0 \in (\theta_Y^-, \theta_Y^+) \mid X \in F_{\theta_0}) = 1 - \gamma$$

Поэтому доверительные интервалы можно использовать для проверки гипотез

Но почему в схеме $\begin{cases} H_0 : a = \bar{x} \\ H_1 : a \neq \bar{x} \end{cases}$ основная гипотеза всегда верна, тогда как выборочно среднее на практике почти всегда не равняется матожиданию. Потому что ...

А вот нефиг такие гипотезы вообще выдвигать

© Блаженов А. В.

Критерий вероятности появления события

$\square P(A) = p$ - вероятность успеха при одном испытании. При достаточно большом количестве испытаний n событие A появилось m раз. Проверяется $H_0 : p = p_0$ против $H_1 : p \neq p_0$
В качестве статистики критерия возьмем величину $K = \frac{m - np_0}{\sqrt{np_0q_0}}$

1. Если H_0 верна, то $K = \frac{m - np}{\sqrt{npq}} \Rightarrow N(0, 1)$ по ЦПТ
2. Lab.

Из тех же соображений $t_{кр}$ - квантиль $N(0, 1)$ уровня $1 - \frac{\alpha}{2}$

$$\begin{cases} H_0 : p = p_0, & \text{если } |K| < t_{кр} \\ H_1 : p \neq p_0 & \text{если } |K| \geq t_{кр} \end{cases}$$

Ех. При посеве 4000 семян 970 всходов оказались рецессивного цвета, а 3030 - доминантного. Проверим гипотезу $H_0 : p = \frac{1}{4}$ - Мендель прав, против $H_1 : p \neq \frac{1}{4}$ - Мендель не прав, для уровня значимости - 0.05

$$K = \frac{m - np_0}{\sqrt{np_0q_0}} = \frac{970 - 4000 \cdot \frac{1}{4}}{\sqrt{4000 \cdot \frac{1}{4} \cdot \frac{3}{4}}} \approx -1.095$$

Так как $|K| = 1.095 < 1.96 = t_{кр}$, то $H_0 : p = \frac{1}{4}$ верна

Лекция 7.

Критерии для проверки гипотез о распределении

Простая параметрическая гипотеза

Пусть имеется выборка (X_1, \dots, X_n) объема n из неизвестного распределения \mathcal{F} . Проверяется простая гипотеза $H_0 : \mathcal{F} = \mathcal{F}_1$ против $H_1 : \mathcal{F} \neq \mathcal{F}_1$, где \mathcal{F}_1 - распределение известного типа с известными нами параметрами $\theta = (\theta_1, \dots, \theta_m)$

I. Критерий Колмогорова

Если \mathcal{F}_1 - абсолютно непрерывное распределение с функцией распределения $F(x)$, то применим критерий

$\square K = \sqrt{n} \sup_x |F^*(x) - F(x)|$, где $F^*(x)$ - выборочная функция распределения

То есть используем теорему Колмогорова: если $H_0 : \mathcal{F} = \mathcal{F}_1$, то $K = \sqrt{n} \sup_x |F^*(x) - F(x)| \Rightarrow \mathcal{K}$ - распределение Колмогорова с функцией распределения $F_{\mathcal{K}}(x) = \sum_{j=-\infty}^{\infty} (-1)^j e^{-2j^2 x^2}$

Для уровня значимости α находим квантиль t_α такой, что $P(\xi \geq t_\alpha) = \alpha$, где $\xi \in \mathcal{K}$

$$\begin{cases} H_0 : \mathcal{F} = \mathcal{F}_1, & \text{если } K < t_\alpha \\ H_1 : \mathcal{F} \neq \mathcal{F}_1, & \text{если } K \geq t_\alpha \end{cases}$$

II. Критерий «хи-квадрат» Пирсона

Пусть выборка разбита на k интервалов A_1, A_2, \dots, A_k , $A_i = [a_{i-1}, a_1)$

n_i - соответствующая частота интервала

При распределении \mathcal{F}_1 теоретические вероятности попадания в эти интервалы $p_i = F_{\mathcal{F}_1}(a_i) - F_{\mathcal{F}_1}(a_{i-1})$. Тогда $n'_i = p_i \cdot n$ - теоретические частоты

В качестве статистики критерия выберем $\chi^2_{\text{набл.}} = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i}$

Th. Пирсона. Если $H_0 : \mathcal{F} = \mathcal{F}_1$ верна, то $\chi^2_{\text{набл.}} \Rightarrow \chi^2_{k-1}$ - распределение «хи-квадрат» с $k - 1$ степенями свободы

Критерий: $\square t_\alpha$ - квантиль χ^2_{k-1} уровня α

$$\begin{cases} H_0 : \mathcal{F} = \mathcal{F}_1, & \text{если } \chi^2_{\text{набл.}} < t_\alpha \\ H_1 : \mathcal{F} \neq \mathcal{F}_1, & \text{если } \chi^2_{\text{набл.}} \geq t_\alpha \end{cases}$$

Nota. Часто обозначают $t_\alpha = \chi^2_{\text{теор.}}$

Nota. При этом частота каждого интервала должна быть не меньше 5, а объем выборки - не меньше 50. Число интервалов лучше брать по формуле Стерджесса

Сложная параметрическая гипотеза

Здесь мы будем проверять гипотезу $H_0 : \mathcal{F} \in \mathcal{F}_\theta$ против $H_1 : \mathcal{F} \notin \mathcal{F}_\theta$, где \mathcal{F}_θ - распределение известного типа с неизвестными параметрами

III. Критерий «хи-квадрат» Фишера

Пусть выборка разбита на k интервалов A_1, A_2, \dots, A_k , $A_i = [a_{i-1}, a_1)$

n_i - соответствующая частота интервала A_i

Def. Оценка максимального правдоподобия по частотам называется значения неизвест-

ных параметров, при которых вероятность появления таких частот является максимальной

Пусть $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_n)$ - оценка максимального правдоподобия по частотам неизвестных параметров. Тогда теоретические вероятности попадания в интервал считаем по формуле $p_i = F_{\mathcal{F}_\theta}(a_i) - F_{\mathcal{F}_\theta}(a_{i-1})$, теоретическая частота - $n'_i = np_i$

В качестве статистики критерия берется функция $\chi^2_{\text{набл.}} = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i}$

Th. Фишера.

Если $H_0 : \mathcal{F} \in \mathcal{F}_\theta$ верна, то $\chi^2_{\text{набл.}} = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i} \Rightarrow \chi^2_{k-m-1}$ - распределение «хи-квадрат», где m - число параметров неизвестного распределения

$\square t_\alpha$ - квантиль χ^2_{k-m-1} уровня α

$$\begin{cases} H_0 : \mathcal{F} = \mathcal{F}_1, & \text{если } \chi^2_{\text{набл.}} < t_\alpha \\ H_1 : \mathcal{F} \neq \mathcal{F}_1, & \text{если } \chi^2_{\text{набл.}} \geq t_\alpha \end{cases}$$

Nota. Часто в качестве оценки неизвестных параметров берется просто оценка максимального правдоподобия

Ex. Имеется выборка в виде вариационного ряда (5.2; ...; 22.8), $n = 120$, при разбиении на $k = 8$ интервалов получили

A_i	[5.2, 7.4)	[7.4, 9.6)	[9.6, 11.8)	[11.8, 14)	[14, 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

Проверить гипотезу о равномерном распределении $H_0 : \mathcal{F} \in U(a, b)$ против $H_1 : \mathcal{F} \notin U(a, b)$ при уровне значимости $\alpha = 0.05$

Дадим оценку параметров методом максимального правдоподобия: $\hat{a} = 5.2$ $\hat{b} = 22.8$

Теоретическая вероятность будет $p'_i = \frac{1}{8}$, теоретическая частота - $n'_i = 15$

$$\chi^2_{\text{набл.}} = \sum_{i=1}^k \frac{(n_i - n'_i)^2}{n'_i} = \frac{(12 - 15)^2}{15} + \frac{(17 - 15)^2}{15} + \frac{(14 - 15)^2}{15} + \frac{(13 - 15)^2}{15} + \frac{(18 - 15)^2}{15} + \frac{(14 - 15)^2}{15} + \frac{(19 - 15)^2}{15} = 3.2$$

При $\alpha = 0.05$ и $S = k - m - 1 = 5$ квантиль χ^2_S уровня α равен $t_\alpha = 11.07$

Так как $\chi^2_{\text{набл.}} < t_\alpha$, нулевая гипотеза о равномерном распределении принимается

Критерии для проверки однородности

IV. Критерий Колмогорова-Смирнова

Пусть имеются 2 независимых выборки (X_1, \dots, X_n) и (Y_1, \dots, Y_m) объемов n и m из неизвестных непрерывных распределений \mathcal{F} и \mathcal{J}

Проверяется $H_0 : \mathcal{F} = \mathcal{J}$ (данные однородны) против $H_1 : \mathcal{F} \neq \mathcal{J}$

В качестве статистики критерия берется функция $K = \sqrt{\frac{nm}{n+m}} \sup_x |F^*(x) - G^*(x)|$, где F^* и G^* - соответствующие выборочные функции распределения

Th. Колмогорова-Смирнова.

Если $H_0 : \mathcal{F} = \mathcal{J}$ верна, то $K \Rightarrow \mathcal{K}$ - распределение Колмогорова

Критерий: t_α - квантиль \mathcal{K} уровня значимости α

$$\begin{cases} H_0 : \mathcal{F} = \mathcal{J}, & \text{если } K < t_\alpha \\ H_1 : \mathcal{F} \neq \mathcal{J}, & \text{если } K \geq t_\alpha \end{cases}$$

Проверки однородности выборок из нормальных совокупностей

V. Критерий Фишера

Пусть имеют две независимые выборки (X_1, \dots, X_n) и (Y_1, \dots, Y_m) объемов n и m из нормальных распределений $X \in N(a_1, \sigma_1^2)$ и $Y \in N(a_2, \sigma_2^2)$

Проверяется $H_0 : \sigma_1 = \sigma_2$ против $H_1 : \sigma_1 \neq \sigma_2$

В качестве статистики критерия берется функция $K = \frac{S_X^2}{S_Y^2}$, где S_X^2, S_Y^2 - соответствующие исправленные дисперсии, причем $S_X^2 \geq S_Y^2$

Th. Если $H_0 : \sigma_1 = \sigma_2$ верна, то $K = \frac{S_X^2}{S_Y^2} \in F(n-1, m-1)$ - распределение Фишера-Снедекера

По пункту 3 основной теоремы $\frac{(n-1)S^2}{\sigma^2} \in \chi_{n-1}^2$. Если H_0 верна, то $K = \frac{S_X^2}{S_Y^2} =$

$$\frac{(n-1)S_X^2 \sigma_2^2}{\sigma_1^2 (m-1) S_Y^2} \frac{(m-1)}{(n-1)} = \frac{\frac{\chi_{n-1}^2}{n-1}}{\frac{\chi_{m-1}^2}{m-1}} \in F(n-1, m-1)$$

t_α - квантиль $F(n-1, m-1)$ уровня α

$$\begin{cases} H_0 : \sigma_1 = \sigma_2, & \text{если } K < t_\alpha \\ H_1 : \sigma_1 \neq \sigma_2, & \text{если } K \geq t_\alpha \end{cases}$$

Nota. Здесь критерий согласия работает чуть иным образом: при верной альтернативной гипотезе $K = \frac{S_X^2}{S_Y^2} \xrightarrow{p} \frac{\sigma_1^2}{\sigma_2^2} > 1$

Если нулевая гипотеза отклоняется, то отклоняется общая гипотеза об однородности. А

если основная гипотеза принимается, то применяем критерий Стьюдента

VI. Критерий Стьюдента

Пусть (X_1, \dots, X_n) и (Y_1, \dots, Y_m) из нормальных распределений $X \in N(a_1, \sigma^2)$ и $Y \in N(a_2, \sigma^2)$

Проверяется $H_0 : a_1 = a_2$ против $H_1 : a_1 \neq a_2$

$$\text{Th. } \sqrt{\frac{nm}{n+m}} \frac{(\bar{x} - a_1) - (\bar{y} - a_2)}{\sqrt{\frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}}} \in T_{n+m-2}$$

По пункту 5 основной теоремы считаем, что числитель и знаменатель независимы

$$\sqrt{\frac{nm}{n+m}} \frac{(\bar{x} - a_1) - (\bar{y} - a_2)}{\sqrt{\frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}}} = \sqrt{\frac{nm}{n+m}} \frac{\frac{\bar{x} - a_1}{\sigma} - \frac{\bar{y} - a_2}{\sigma}}{\sqrt{\frac{(n-1)S_X^2 + (m-1)S_Y^2}{\sigma^2(n+m-2)}}}$$

По пункту 1 основной теоремы $\sqrt{n} \frac{\bar{x} - a_1}{\sigma}, \sqrt{m} \frac{\bar{y} - a_2}{\sigma} \in N(0, 1) \Rightarrow$

$$\frac{\bar{x} - a_1}{\sigma} \in N\left(0, \frac{1}{n}\right), \frac{\bar{y} - a_2}{\sigma} \in N\left(0, \frac{1}{m}\right) \Rightarrow$$

$$\frac{\bar{x} - a_1}{\sigma} - \frac{\bar{y} - a_2}{\sigma} \in N\left(0, \sqrt{\frac{n+m}{nm}}\right) \Rightarrow$$

$$\sqrt{\frac{nm}{n+m}} \left(\frac{\bar{x} - a_1}{\sigma} - \frac{\bar{y} - a_2}{\sigma} \right) \in N(0, 1)$$

По пункту 3 основной теоремы $\frac{(n-1)S_X^2}{\sigma^2} \in \chi_{n-1}^2, \frac{(m-1)S_Y^2}{\sigma^2} \in \chi_{m-1}^2 \Rightarrow$

$$\frac{(n-1)S_X^2 + (m-1)S_Y^2}{\sigma^2(n+m-2)} \in \frac{\chi_{n+m-2}^2}{m+n-2}$$

$$\text{Из этого } \sqrt{\frac{nm}{n+m}} \frac{(\bar{x} - a_1) - (\bar{y} - a_2)}{\sqrt{\frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}}} \in \frac{N(0, 1)}{\frac{\chi_{n+m-2}^2}{n+m-2}} = T_{n+m-2}$$

В качестве статистики возьмем $\sqrt{\frac{nm}{n+m}} \frac{(\bar{x} - a_1) - (\bar{y} - a_2)}{\sqrt{\frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}}}$, по теореме при $a_1 = a_2$ получаем,

что $K \in T_{n+m-2}$

Если верна альтернативная гипотеза, то $K \rightarrow \infty$

Критерий: t_α - квантиль $|T_{n+m-2}|$ уровня α

$$\begin{cases} H_0 : a_1 = a_2, & \text{если } K < t_\alpha \\ H_1 : a_1 \neq a_2, & \text{если } K \geq t_\alpha \end{cases}$$

Nota. Если при обоих критериях согласились с нулевой гипотезой, то соглашаемся с гипотезой об однородности выборок

Nota. Критерий хорошо работает, если выборки из нормальных распределений (или очень близких к ним)

Лекция 8.

Статистическая зависимость

Def. Зависимость называется статистической, если изменение одной случайной величины вызывает изменение распределения другой. Если при этом изменяется среднее значение другой случайной величины, то такая зависимость называется корреляционной. Если при увеличении одной случайной величины среднее значение другой также увеличивается, то говорят, что имеет место прямая корреляция. Аналогично, если уменьшается, то – обратная

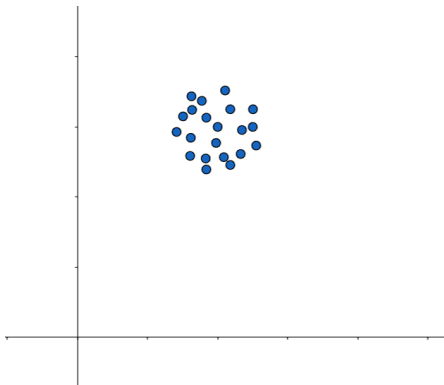
Корреляционное облако

Пусть в ходе n экспериментов появились значения двух случайных величин X и Y :

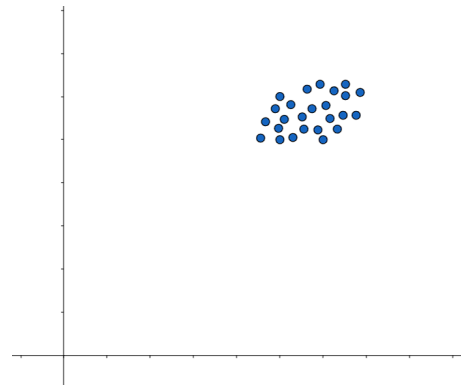
$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

Нанеся точки на координатную плоскость, получаем корреляционное облако, о виде которого можно делать предположения о наличии/отсутствии связи

Ex.



Здесь, возможно, нет зависимости



Здесь можно предположить прямую корреляцию

Корреляционная таблица

Пусть даны данные X и Y при n экспериментов. Эти данные удобно представить в виде корреляционной таблицы: по вертикали отмечают различные значения x , а по горизонтали – y , в клетках таблицы отмечаются частота появления n_{xy}

Ex. $n = 50$

$X \backslash Y$	10	20	30	40	n_x	\bar{y}_x
2	7	3	0	0	10	13
4	3	10	10	2	25	4.4
6	0	2	10	3	15	30.67
n_y	10	15	20	5	$\Sigma 50$	

По диагонали таблицы можно предположить, что корреляция есть

Имеет смысл вычислить условное среднее по формуле $\bar{y}(x) = \frac{1}{n_x} \sum n_{xy} y_i$. Так как в нашем примере условные средние растут с ростом x , то имеет место прямая корреляция

Nota. Если данных много или X и Y – непрерывные случайные величины, то лучше составить интервальную корреляционную таблицу: разбить случайные величины на интервалы, по вертикали отметить интервалы $[a_{i-1}, a_i)$ случайной величины X , по горизонтали – $[b_{j-1}, b_j)$ случайной величины Y , в клетках отметить частоты $n_{ij} : [a_{i-1}, a_i) \times [b_{j-1}, b_j)$. В дальнейшем интервалы можно заменить их серединами

Критерий «хи-квадрат» для проверки независимости

Пусть выборка $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ представлена в виде интервальной корреляционной таблицы. Случайная величина X разбита на k интервалов, а Y – на m интервалов

Обозначим v_i – частота i -ого интервала $[a_{i-1}, a_i)$ случайной величины X , v_j – частота j -ого интервала $[b_{j-1}, b_j)$ случайной величины Y , v_{ij} – число точек в $[a_{i-1}, a_i) \times [b_{j-1}, b_j)$

$X \backslash Y$	$[b_0, b_1)$	$[b_1, b_2)$	\dots	$[b_{m-1}, b_m)$	v_i
$[a_0, a_1)$	v_{11}	v_{12}	\dots	v_{1m}	$v_{1\cdot}$
\dots					
$[a_{k-1}, a_k)$	v_{k1}	v_{k2}	\dots	v_{km}	$v_{k\cdot}$
$v_{\cdot j}$	10	15	20	5	Σn

Проверяется основная гипотеза $H_0 : X$ и Y независимы против $H_1 = \overline{H_0} : X$ и Y зависимы

Если H_0 верна, то $p_{ij} = P(X \in [a_{i-1}, a_i), Y \in [b_{j-1}, b_j)) = P(X \in [a_{i-1}, a_i)) \cdot P(Y \in [b_{j-1}, b_j))$

Тогда по закону больших чисел $\frac{v_{i\cdot}}{n} \xrightarrow{P} p_{i\cdot}$, $\frac{v_{\cdot j}}{n} \xrightarrow{P} p_{\cdot j}$

Поэтому основанием для отклонения основной гипотезы будет заметная разница между величинами $\frac{v_{i\cdot} \cdot v_{\cdot j}}{n}$ и $\frac{v_{ij}}{n}$ или v_{ij} и $\frac{1}{n} v_{i\cdot} v_{\cdot j}$

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

Th. Если H_0 верна, то $K \Rightarrow H_{(k-1)(m-1)}$

Пусть t_α – квантиль $H_{(k-1)(m-1)}$ уровня α , тогда

$$\begin{cases} H_0 : X \text{ и } Y \text{ независимы, если } K < t_\alpha \\ H_0 : X \text{ и } Y \text{ зависимы, если } K \geq t_\alpha \end{cases}$$

Nota. Для работы критерия необходимо, чтобы частота в каждой клетке была больше 5, а объем выборки был достаточно большой

Однофакторный дисперсионный анализ

Предположим, что на случайную величину X (результат) может влиять фактор Z (необязательно, что Z – случайная величина, эксперимент может быть управляемым)

Пусть при различных « k уровней» фактора Z получено k независимых выборок случайной величины X : $X^{(1)} = (X_1^{(1)}, \dots, X_{n_1}^{(1)}), \dots, X^{(k)} = (X_1^{(k)}, \dots, X_{n_k}^{(k)})$

Всего было получено $n = \sum_{i=1}^k n_i$ значений

Nota. В общем говоря, распределение этих выборок отличается, поэтому эти выборки разных случайных величин

Общая, внутригрупповая и межгрупповая дисперсии

Для каждой выборки вычислим выборочное среднее и дисперсию: $\bar{x}^{(j)} = \frac{1}{n_j} \sum_{i=1}^{n_j} X_i^{(j)}, D^{(j)} =$

$$\frac{1}{n_j} \sum_{i=1}^{n_j} (X_i^{(j)} - \bar{x}^{(j)})^2$$

Объединим все выборки в общую и также вычислим выборочное среднее и дисперсию:

$$\bar{x} = \frac{1}{n} \sum_{i,j} x_i^{(j)} = \frac{1}{n} \sum_{j=1}^k n_j \cdot \bar{x}^{(j)} - \text{общее среднее}$$

$$D_O = \frac{1}{n} \sum_{i,j} (X_i^{(j)} - \bar{x})^2 - \text{общая дисперсия}$$

Def. Внутригрупповой (или остаточной) дисперсией называется среднее групповых дисперсий:

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

Def. Межгрупповой (или факторной) дисперсией называется величина $D_M = \frac{1}{n} \sum_{j=1}^k n_j (\bar{x} - \bar{x}^{(j)})^2$

Th. О разложении дисперсии. Общая дисперсия равна сумме внутригрупповой и межгрупповой дисперсией: $D_O = D_B + D_M$

Смысл: внутригрупповая дисперсия показывает средний (случайный) разброс внутри выборок, межгрупповая - насколько отличаются среднее при различных уровнях фактора, то есть именно эта величина отражает влияния фактора

Вывод по наличию корреляции можно сделать, если доля D_M достаточно велика

Проверка гипотезы о влиянии фактора

Предполагаем, что X имеет нормальное распределение и фактор Z может влиять только на ее математическое ожидание, но на дисперсию и тип распределения, поэтому можно считать, что данные независимых k выборок при различных уровнях фактора Z также имеют нормальное распределение с одинаковой дисперсией: $X^{(j)} \in N(a_j, \sigma^2)$

Проверяется основная гипотеза $H_0 : a_1 = a_2 = \dots = a_k$ (фактор не оказывает влияния) против $H_1 = \overline{H_0}$: есть влияние

По пункту 3 основной теоремы $\sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 = \frac{nD^*}{\sigma^2} \in H_{n-1}$

Из этого $\frac{n_j D^{(j)}}{\sigma^2} \in H_{n_j-1} \quad \forall 1 \leq j \leq k$

Так как распределение «хи-квадрат» устойчиво относительно суммирования, то $\sum_{j=1}^k \frac{n_j D^{(j)}}{\sigma^2} \in H_{n-k}$, так как $(n_1 - 1) + \dots + (n_k - 1) = n - k$

Пусть основная гипотеза верна, тогда все данные можно считать выборкой одной случайной величины и по пункту 3 $\frac{nD_O}{\sigma^2} \in H_{n-1}$

Согласно теореме о разложении дисперсии $D_O = D_B + D_M$, тогда $\frac{nD_O}{\sigma^2} = \frac{nD_B}{\sigma^2} + \frac{nD_M}{\sigma^2}$

Так как $\frac{nD_O}{\sigma^2} \in H_{n-1}$, $\frac{nD_B}{\sigma^2} \in H_{n-k}$, то $\frac{nD_M}{\sigma^2} \in H_{k-1}$

Тогда при верной основной гипотезе получим, что $\frac{nD_M}{\sigma^2} \frac{\sigma^2(n-k)}{nD_B} = \frac{n-k}{k-1} \frac{D_M}{D_B} \in F(k-1, n-k)$ – распределение Фишера-Снедекера со степенями $k-1$ и $n-k$

В качестве статистики берется $K = \frac{n-k}{k-1} \frac{D_M}{D_B}$, в качестве критической точки t_α – квантиль $F(k-1, n-k)$ уровня α

$\begin{cases} H_0 : a_1 = a_2 = \dots = a_k \text{ (фактор оказывает влияние), если } K < t_\alpha \\ H_1 : \text{фактор влияния не оказывает, если } K \geq t_\alpha \end{cases}$

Лекция 9.

Исследование статистической корреляции

Математическая модель регрессии

Пусть случайная величина X зависит от случайной величины Z (необязательно случайной)

Def. Регрессией X на Z называется функция $f(z) = E(X|Z = z)$. Она показывает зависимость среднего значения X от значения Z

Уравнение $x = f(z)$ называется уравнением регрессии, а график этой функции - линия регрессии

Пусть при n экспериментах при значениях Z_1, Z_2, \dots, Z_n фактора Z наблюдались значения X_1, X_2, \dots, X_n случайной величины X

Обозначим через ε_i разницу между экспериментальным и теоретическими значениями случайной величины X , то есть $\varepsilon_i = X_i - f(z_i)$

ε - это случайный член модели или так называемая теоретическая ошибка

Nota. Обычно можно считать, что ε_i независимы друг от друга и имеет нормальное распределение с $a = 0$, так как $E\varepsilon_i = E(X_i - f(Z_i)) = E(X|Z = Z_i) - E(X|Z = Z_i) = 0$

Цель: нам нужно по экспериментальным данным $(z_1, x_1), \dots, (z_n, x_n)$ как можно лучше оценить функцию $f(z)$

Nota. При этом предполагая (часто из теории), что $f(z)$ - функция определенного вида, но параметры которой неизвестны. Если нет, то начинаем подбирать модели самого простого вида. В противном случае, наилучшим решением была бы кривая, проходящая через все точки

Метод наименьших квадратов

Пусть известен из теории вид функции $f(z)$. Метод наименьших квадратов состоит в выборе параметров $f(z)$ таким образом, чтобы минимизировать сумму квадратов ошибок $\sum_{i=1}^n \varepsilon_i^2 =$

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

Def. Пусть θ - набор неизвестных параметров функции $f(z)$. Оценка $\hat{\theta}$ параметра θ , при которой достигается минимум $\sum_{i=1}^n \varepsilon_i^2$, называется оценкой метода наименьших квадратов (или ОМНК)

Линейная парная регрессия

Пусть имеется теоретическая модель линейной регрессии

$f(z) = \alpha + \beta z + \varepsilon$ - теоретическая модель, где ε - теоретическая ошибка отражающая влияние невключенных в модель факторов, возможной нелинейности, ошибок измерения и просто случая

Пусть $(z_1, x_1), \dots, (z_n, x_n)$ - экспериментальные данные. По ним методом наименьших квадратов строим экспериментальную модель линейной регрессии $f(z) = a + bz$, где a и b - ОМНК параметров α и β

$\hat{\varepsilon}_i = X_i - f(Z_i) = X_i - (a + bZ_i)$ - экспериментальная ошибка

Найдем ОМНК параметров α и β

$$\begin{aligned}
\sum_{i=1}^n \hat{\varepsilon}_i^2 &= \sum_{i=1}^n (X_i - (a + bZ_i))^2 \\
\frac{\partial}{\partial a} \sum_{i=1}^n \hat{\varepsilon}_i^2 &= \sum_{i=1}^n -2(X_i - a - bZ_i) = -2 \sum_{i=1}^n X_i + 2 \sum_{i=1}^n a + 2b \sum_{i=1}^n Z_i = 2(n\bar{x} - na - bn\bar{z}) \\
\frac{\partial}{\partial b} \sum_{i=1}^n \hat{\varepsilon}_i^2 &= \sum_{i=1}^n -2Z_i(X_i - a - bZ_i) = -2 \sum_{i=1}^n X_i Z_i + 2 \sum_{i=1}^n a Z_i + 2b \sum_{i=1}^n Z_i^2 = 2(n\bar{x}\bar{z} - an\bar{z} - bn\bar{z}^3) \\
\begin{cases} -2(n\bar{x} - na - bn\bar{z}) = 0 \\ -2(n\bar{x}\bar{z} - an\bar{z} - bn\bar{z}^3) = 0 \end{cases} &\iff \begin{cases} a + b\bar{z} = \bar{x} \\ a\bar{z} + b\bar{z}^2 = \bar{x}\bar{z} \end{cases}
\end{aligned}$$

Получили систему линейных уравнений. Будем называть ее нормальной системой. При решении получаем:

$$\begin{cases} a = \bar{x} - b\bar{z} \\ (\bar{x}b\bar{z})\bar{z} + b\bar{z}^2 = \bar{x}\bar{z} \end{cases} \iff \begin{cases} a = \bar{x} - b\bar{z} \\ b = \frac{\bar{x}\bar{z} - \bar{x}\bar{z}}{\bar{z}^2} \end{cases} \quad - \text{ОМНК}$$

Запишем уравнение линейной регрессии в удобном виде: $\bar{x}_z = f(z) = E(X|Z = z)$

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

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

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

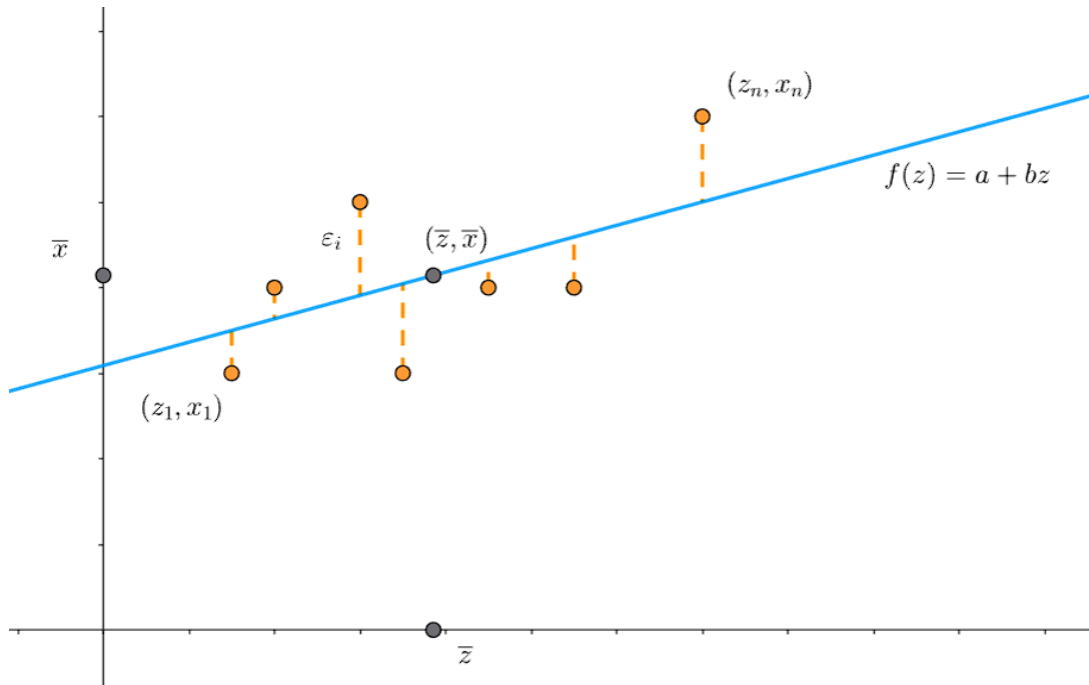
$$\bar{x}_z - \bar{x} = \frac{\hat{\sigma}_x}{\hat{\sigma}_z} \frac{\bar{z}\bar{x} - \bar{x}\bar{z}}{\hat{\sigma}_z \hat{\sigma}_x} (z - \bar{z}) = \frac{\hat{\sigma}_x}{\hat{\sigma}_z} \hat{r} (z - \bar{z}), \text{ где } \hat{r} - \text{выборочный коэффициент линейной корреляции}$$

Или $\frac{\bar{x}_z - \bar{x}}{\hat{\sigma}_x} = \hat{r} \frac{z - \bar{z}}{\hat{\sigma}_z}$ - выборочное уравнение линейной регрессии

Nota. Прямая регрессии проходит через точку из выборочных средних

Nota. При $n \rightarrow \infty$ $\bar{x} \rightarrow EX, \bar{z} \rightarrow EZ, \hat{\sigma}_x \rightarrow \sigma_x, \hat{\sigma}_z \rightarrow \sigma_z, \bar{x}_z \rightarrow E(X|Z = z), \hat{r} \rightarrow r$, получаем $\frac{E(X|Z = z) - EX}{\sigma_x} = r \frac{z - EZ}{\sigma_z}$ - теоретическое уравнение линейной регрессии

Геометрический смысл линии регрессии



Суть МНК: находим такую прямую, чтобы сумма квадратов длин этих отрезков (по сути отклонений) была минимальна (или дисперсия экспериментальных данных относительно прямой была минимальна)

Выборочный коэффициент линейной корреляции

Def. $\hat{r} = \frac{\overline{zx} - \bar{x}\bar{z}}{\hat{\sigma}_z \hat{\sigma}_x}$ называется выборочным коэффициентом линейной корреляции. Ясно, что она будет точечной оценкой теоретического коэффициента линейной корреляции. Также \hat{r} является несмещенной оценкой

Поэтому выборочный коэффициент корреляции характеризует силу линейной связи. Знак коэффициента показывает направления корреляции (прямая или обратная)

Силу связи можно примерно оценить от шкале Чеддока:

Количественная мера \hat{r}	Качественная мера
0.1 – 0.3	Слабая
0.3 – 0.5	Умеренная
0.5 – 0.7	Заметная
0.7 – 0.9	Высокая
> 0.9	Весьма высокая

Проверка гипотезы о значимости выборочного коэффициента корреляции

Пусть (Z, X) распределена нормально. По выборке объема n вычислен выборочный коэффициент корреляции \hat{r} , а r - теоретический коэффициент корреляции

Проверяется $H_0 : r = 0$ (выборочный коэффициент корреляции статистически незначим) против $H_1 : r \neq 0$ (коэффициент статистически значим)

$$\text{Если } H_0 \text{ верна, то } K = \frac{\hat{r}\sqrt{n-2}}{\sqrt{n-\hat{r}^2}} \in T_{n-2} \text{ - распределение Стьюдента с степенью } n-2$$

Получаем критерий. Пусть t_α - квантиль $|T_{n-2}|$ (двухстороннее распределение Стьюдента) уровня α

$$\begin{cases} H_0 : r = 0, & \text{если } |K| < t_\alpha \\ H_1 : r \neq 0, & \text{если } |K| \geq t_\alpha \end{cases}$$

Надо понимать, что корреляция - более тонкое понятие, чем зависимость

А термин *регрессия* получил свое название чисто исторически: статистик Гальтон в 1886 году исследовал зависимость роста детей от роста родителей

$$E(P_{\text{сына}} | Z_{\text{отца}} = Z_1, Z_{\text{матери}} = Z_1) = 0.27Z_1 + 0.2Z_2 + \text{const}$$

$$E(P_{\text{дочери}} | Z_{\text{отца}} = Z_1, Z_{\text{матери}} = Z_1) = \frac{1}{1.08} P_{\text{сына}}$$

Дальше он заметил, что при у самых высоких родителей рост детей был меньше относительно них (скатывался к среднему, происходил регресс)

Позже исследовали экономические результаты фирм, показатели спортсменов, которые после успешного сезона уменьшались, после чего появлялось куча теорий. Сейчас все это объясняется простым случаем

Выборочное корреляционное отношение

Выборочный коэффициент корреляции характеризует только силу линейной связи. Следующий подход основан на однофакторном дисперсионном анализе

Пусть есть k выборок случайной величины X при k различных уровнях фактора Z . Вычислены общая, внутригрупповая и межгрупповая дисперсии. По теореме $D_O = D_M + D_B$

Def. Выборочным корреляционным отношением X на Z называется величина $\eta_{X,Z} = \sqrt{\frac{D_M}{D_O}}$

Свойства:

1. $0 \leq \eta_{X,Z} \leq 1$ ($D_M, D_O \geq 0$)
2. Если $\eta = 1$, то $D_M = D_O \implies D_B = 0$, имеем функциональную зависимость X от Z
3. Если $\eta = 0$, то $D_M = 0 \implies$ корреляция отсутствует

4. $\eta \geq |\hat{r}|$
5. Если $\eta = |\hat{r}|$, то все точки экспериментальных данных лежат на прямой линейной регрессии (то есть данная линейная модель является идеальной)

Лекция 10.

Свойство ковариации

Мет. Ковариацией случайных величин X и Y называется величина $\text{cov}(X, Y) = E((X - EX)(Y - EY)) = E(XY) - EX \cdot EY$

Ковариация является индикатором наличия направления связи между двумя случайными величинами

Пусть имеется $(X_1, Y_1), \dots, (X_n, Y_n)$ случайных величин X и Y

Def. Выборочной ковариацией называется величина $\widehat{\text{cov}}(X, Y) = \overline{xy} - \bar{x} \cdot \bar{y}$

По Закону Больших Чисел ясно, что $\widehat{\text{cov}}(X, Y) \rightarrow \text{cov}(X, Y)$, поэтому выборочная ковариация является оценкой

Th. Выборочная ковариация является точечной состоятельной, но смещенной оценкой ковариации. Несмещенной оценкой будет $\frac{n}{n-1} \widehat{\text{cov}}(X, Y)$

Ковариация и выборочная ковариация обладают свойствами

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. $\text{cov}(X, X) = D(X)$, $\widehat{\text{cov}}(X, X) = D^*(X)$
6. $D(X + Y) = DX + DY + 2 \text{cov}(X, Y)$

Nota. В дальнейшем под $\text{cov}(X, Y)$ будет пониматься выборочная ковариация

Анализ модели линейной парной регрессии

Пусть при n экспериментах получены значения случайных величин X и Z : $(X_1, Z_1), \dots, (X_n, Z_n)$

Пусть $X = \alpha + \beta Z + \varepsilon$ - теоретическая модель линейной регрессии, где ε - случайная величина, отражающая влияние не включенных факторов, нелинейность модели, ошибок измерений и просто случая.

Пусть построили с помощью метода наименьших квадратов выборочное уравнение линейной регрессии $\hat{X} = a + bZ$

Обозначим $\hat{\varepsilon}_i = X_i - \hat{X}_i$ - экспериментальная ошибка, разница между наблюдаемыми значениями и вычисляемыми по модели

Тогда $X_i = \hat{X}_i + \hat{\varepsilon}_i$ или $X_i = a + bZ_i + \hat{\varepsilon}_i$, где a и b - точечные оценки параметров α и β

Свойства $\hat{\varepsilon}_i$:

1. $\overline{\hat{\varepsilon}_i} = 0$

$$a = \bar{X} - b\bar{Z} \implies a + b\bar{Z} = \bar{X} \implies \overline{\hat{\varepsilon}_i} = \overline{X_i - (a + bZ_i)} = \bar{X} - \overline{a + bZ_i} = \bar{X} - \bar{X} = 0$$

2. $\text{cov}(\hat{X}, \hat{\varepsilon}) = 0$

$$\begin{aligned} b &= \overline{xz} - \bar{x} \cdot \bar{z} \hat{\sigma}_z^2 = \overline{\text{cov}(X, Z)} D(Z) \implies \text{cov}(X, Z) - bD(Z) = 0 \\ \text{cov}(\hat{X}, \hat{\varepsilon}) &= \text{cov}(a + bZ, X - a - bZ) = \text{cov}(bZ, X - bZ) = \text{cov}(bZ, x) - \text{cov}(bZ, bZ) = \\ &= b \text{cov}(Z, X) - b^2 D(Z) = b(\text{cov}(Z, X) - bD(Z)) = 0 \end{aligned}$$

Анализ дисперсии результата

Def. $D(X) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ - дисперсия наблюдаемых значений

Def. $\hat{D}(X) = \frac{1}{n} \sum_{i=1}^n (\hat{X}_i - \bar{X})^2$ - дисперсия расчетных значений

Def. $D(\hat{\varepsilon}) = \frac{1}{n} \sum_{i=1}^n (\hat{\varepsilon}_i)^2$ - дисперсия остатков

Так как $X = \hat{X} + \hat{\varepsilon}$, $\text{cov}(\hat{X}, \hat{\varepsilon}) = 0$, то $D(X) = D(\hat{X}) + D(\hat{\varepsilon}) + 2 \text{cov}(\hat{X}, \hat{\varepsilon}) = D(\hat{X}) + D(\hat{\varepsilon})$

$$\text{Th. } D(X) = D(\hat{X}) + D(\hat{\varepsilon})$$

Очевидно, что качество модели будет тем лучше, чем меньше будет дисперсия остатков

Def. Коэффициентом детерминации R^2 называется величина $R^2 = \frac{D(\hat{X})}{D(X)}$ или $R^2 = 1 - \frac{D(\hat{\varepsilon})}{D(X)}$

Nota. Смысл R^2 - доля объясненной дисперсии, а $1 - R^2$ - доля необъясненной дисперсии

Свойства:

1. $0 \leq R^2 \leq 1$
2. Если $R^2 = 1$, то $D(\hat{\varepsilon}) = 0 \implies \hat{\varepsilon}_i = \bar{\hat{\varepsilon}}_i = 0$, то есть точки лежат строго на линии регрессии, модель идеальна

3. Если $R^2 = 0$, то $D(\hat{X}) = 0 \implies \hat{X} = \bar{x}$, то есть получаем примитивную, ничего не объясняющую модель

Чем больше R^2 , тем лучше качество модели

Проверка гипотезы о значимости уравнения регрессии

Проверяется $H_0 : R_{\text{теор}}^2 = 0$ (уравнение регрессии статистически не значимо) против $H_1 : R_{\text{теор}}^2 \neq 0$

Th. Если H_0 верна, то $F = \frac{R^2(n-2)}{1-R^2} \in F(1, n-2)$

Пусть t_α - квантиль $F(1, n-2)$ уровня α , тогда:

$$\begin{cases} H_0 : R_{\text{теор}}^2 = 0 & \text{если } F < t_\alpha \\ H_0 : R_{\text{теор}}^2 \neq 0 & \text{если } F \geq t_\alpha \end{cases}$$

Nota. Если $H_0 : R_{\text{теор}}^2 = 0$, то $H_0 : \beta = 0$

Связь между коэффициентом детерминации и коэффициентом линейной корреляции

1. $\sqrt{R^2} = r_{\hat{X}, X}$ - коэффициент корреляции между \hat{X} и X

$$r_{\hat{X}, X} = \frac{\text{cov}(\hat{X}, X)}{\sqrt{D(\hat{X})D(X)}} = \frac{\text{cov}(\hat{X}, \hat{X} + \varepsilon)}{\sqrt{D(\hat{X})D(X)}} = \frac{D(\hat{X}) + \text{cov}(\hat{X}, \varepsilon)}{\sqrt{D(\hat{X})D(X)}} \stackrel{0}{=} \sqrt{\frac{D(\hat{X})}{D(X)}} = R$$

2. $r_{\hat{X}, X} = |r_{X, Z}|$

$$\begin{aligned} \text{cov}(\hat{X}, X) &= \text{cov}(a + bZ, X) = b \text{cov}(Z, X) \\ D(\hat{X}) &= D(a + bZ) = b^2 D(Z) \\ r_{\hat{X}, X} &= \frac{\text{cov}(\hat{X}, X)}{\sqrt{D(\hat{X})D(X)}} = \frac{b \text{cov}(Z, X)}{\sqrt{b^2 D(Z)D(X)}} = \left| \frac{\text{cov}(X, Z)}{\sqrt{D(Z)D(X)}} \right| = |r_{X, Z}| \end{aligned}$$

Следствие 1: в случае линейной парной регрессии коэффициент детерминации равен квадрату коэффициенту корреляции

Следствие 2: в случае линейной парной регрессии совпадают результаты проверки гипотез $H_0 : R_{\text{теор}}^2 = 0 \iff H_0 : r = 0 \iff H_0 : \beta = 0$

Теорема Гаусса-Маркова

Th. Пусть $X_i = \alpha + \beta Z_i + \varepsilon_i$ - теоретическая модель регрессии

$X = a + bZ$ - модель, полученная по методу наименьших квадратов

Если выполнено условия:

- а) Случайные члены ε_i независимые случайные величины, имеющие одинаковое нормальное распределение $\varepsilon_i \in N(0, \sigma^2)$
- б) Случайные величины ε_i и Z_i - независимы

Тогда a и b - состоятельные, несмещенные, эффективные оценки параметров α и β , то есть

- 1. Состоятельность: $a \xrightarrow[n \rightarrow \infty]{p} \alpha, b \xrightarrow[n \rightarrow \infty]{p} \beta$
- 2. Несмещенность: $Ea = \alpha, Eb = \beta$
- 3. Наименьшая дисперсия, равная:

$$Da = \frac{\overline{z^2} \sigma^2}{nD(Z)}, Db = \frac{\sigma^2}{nD(Z)}$$

Nota. Если не выполняется условие а), то есть ошибки зависимы или имеют разную дисперсию, то оценки становятся неэффективными. Если не выполнено условие б), то оценки становятся смещенными и несостоятельными

Стандартные ошибки коэффициентов регрессии

Из теоремы видим, что Da и Db зависят от дисперсии σ^2 случайного члена. По экспериментальным ошибкам получаем оценку данной дисперсии:

$$D(\hat{\varepsilon}) = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 \xrightarrow[n \rightarrow \infty]{p} \sigma^2$$

Однако эта оценка является смещенной:

$$E(D(\hat{\varepsilon})) = \frac{n-2}{n} \sigma^2$$

Поэтому несмещенной оценкой дисперсии σ^2 является величина $S^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2$

Def. Величина S называется стандартной ошибкой регрессии

Смысл: характеризует разброс наблюдаемых значений вокруг линии регрессии

Nota. Заменим в теореме Гаусса-Маркова σ^2 на S^2 , получаем оценки дисперсий Da и Db :

$$S_a^2 = \frac{\overline{z^2} S^2}{nD(z)}, S_b^2 = \frac{S^2}{nD(Z)}$$

Def. S_a и S_b называются стандартными ошибками коэффициентов регрессии

Прогнозирование регрессионных моделей

Пусть $X = \alpha + \beta Z + \varepsilon$ - теоретическая модель

$\hat{X} = a + bZ$ - модель МНК, построенная по выборке объема n

С помощью данной модели надо дать прогноз значения X_p при заданном значении Z_p и оценить качество прогноза

Теоретическое значение - $X_p = \alpha + \beta Z_p + \varepsilon$, а точечный прогноз $\hat{X}_p = a + bZ_p$

Разность между ними $\Delta_p = \hat{X}_p - X_p$ называется ошибкой предсказания

Свойства Δ_p :

1. $E\Delta_p = 0$

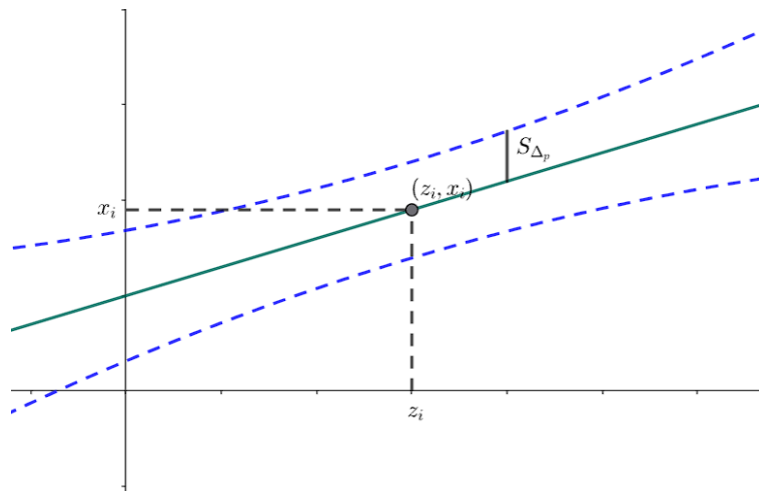
2. $D(\Delta_p) = \left(1 + \frac{1}{n} + \frac{(Z_p - \bar{z})^2}{nDZ}\right) \sigma^2$

Заменяя σ^2 на S^2 , получим стандартную ошибку прогноза: $S_{\Delta_p} = S \sqrt{1 + \frac{1}{n} + \frac{(Z_p - \bar{z})^2}{nDZ}}$

3. $D(\Delta_p) > \sigma^2$ - то есть точность прогноза ограничена случайным членом ε

4. При $n \rightarrow \infty$ $D(\Delta_p) \xrightarrow{p} \sigma^2$ - качество модели тем лучше, чем больше объем выборки

5. Чем больше Z_p отклоняется от \bar{z} , тем хуже качество прогноза. Наилучшее качество в точке $Z_p = \bar{z}$: $D(\Delta_p) = \left(1 + \frac{1}{n}\right) \sigma^2$



Доверительные интервалы прогноза и коэффициентов уравнения линейной регрессии

Пусть t_γ - квантиль $|T_{n-2}|$ уровня γ

Тогда доверительные интервалы надежности γ для параметров α и β :

$$\alpha : (a - t_\gamma S_a; a + t_\gamma S_a)$$

$$\beta : (b - t_\gamma S_b; b + t_\gamma S_b)$$

Доверительный интервал для прогноза X_p : $(\hat{X}_p - t_\gamma S_{\Delta_p}; \hat{X}_p + t_\gamma S_{\Delta_p})$

Лекция 11.

Математическое ожидание и дисперсия случайного вектора

Def. $E\vec{X} = \begin{pmatrix} EX_1 \\ \vdots \\ EX_n \end{pmatrix}$ - математическое ожидание случайного вектора

Def. Дисперсией или матрицей ковариаций называется $D\vec{X} = E((\vec{X} - E\vec{X})(\vec{X} - E\vec{X})^T)$, элементами которой $d_{ij} = \text{cov}(X_i, X_j)$, $d_{ii} = D(X_i)$

Свойства:

1. $E(A\vec{X}) = AE\vec{X}$
2. $E(\vec{X} + \vec{B}) = E\vec{X} + \vec{B}$
3. $D(A\vec{X}) = AD\vec{X}A^T$
4. $D(\vec{X} + \vec{B}) = D\vec{X}$

Уравнение общей регрессии

Пусть результат X зависит от k факторов Z_1, \dots, Z_k . Рассматриваем теоретическую модель линейной регрессии:

$X = \beta_1 Z_1 + \beta_2 Z_2 + \dots + \beta_k Z_k + \varepsilon$, где $\vec{Z} = \begin{pmatrix} Z_1 \\ \vdots \\ Z_k \end{pmatrix}$ - вектор факторов, $\vec{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$ - вектор коэффициентов регрессии

Пусть проведено $n \geq k$ экспериментов, $\vec{Z}^{(i)} = \begin{pmatrix} Z_1^{(i)} \\ \vdots \\ Z_k^{(i)} \end{pmatrix}$ - значения факторов при i -ом эксперименте,

$\vec{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$ - соответствующие значения результатов

Согласно модели:

$$\begin{cases} X_1 = \beta_1 Z_1^{(1)} + \beta_2 Z_2^{(2)} + \dots + \beta_k Z_k^{(1)} + \varepsilon_1 \\ X_1 = \beta_1 Z_1^{(1)} + \beta_2 Z_2^{(2)} + \dots + \beta_k Z_k^{(1)} + \varepsilon_1 \\ \dots \\ X_n = \beta_1 Z_1^{(n)} + \beta_2 Z_2^{(n)} + \dots + \beta_k Z_k^{(n)} + \varepsilon_n \end{cases}$$

Или в матричной форме: $\vec{X} = Z^T \vec{\beta} + \vec{\varepsilon}$, где $Z = \begin{pmatrix} Z_1^{(1)} & Z_1^{(2)} & \dots & Z_1^{(n)} \\ Z_2^{(1)} & Z_2^{(2)} & \dots & Z_2^{(n)} \\ \vdots & \vdots & \ddots & \vdots \\ Z_k^{(1)} & Z_k^{(2)} & \dots & Z_k^{(n)} \end{pmatrix}$ - матрица плана, $\vec{\varepsilon}$ - вектор

теоретических ошибок

Наша цель такова: по данной матрице плана Z и вектора результатов \vec{X} дать оценки неизвестных параметров регрессии β_i и параметров распределения ошибки ε

Nota. Заметим, что у данной модели мы не теряем свободный член b_0 , так как при необходимости можно считать, что первый фактор тождественен единицы. То есть первая строка матрицы плана будет состоять из единиц

Метод наименьших квадратов и нормальные уравнения

Будем считать, что выполнено условие Cond.1, что ранг матрицы $\text{rang } Z = k$, то есть все строки матрицы плана независимы

Введем матрицу $A = ZZ^T$. Ее свойства:

1. A - квадратная и симметричная
2. A - положительно определенная
3. $\exists B = \sqrt{A}$, то есть $B^2 = A$

Пусть $\vec{B} = \begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix}$ - вектор оценок $\vec{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$

Тогда эмпирическая модель регрессии $\vec{X} = Z^T \vec{B}$, $\varepsilon_i = X_i - \hat{X}_i$ - экспериментальная ошибка, или $\vec{\varepsilon} = \begin{pmatrix} \hat{\varepsilon}_1 \\ \vdots \\ \hat{\varepsilon}_n \end{pmatrix} = \vec{X} - Z^T \vec{B}$ - вектор экспериментальных ошибок

По методу наименьших квадратов подбираем \vec{B} таким образом, чтобы $L(\vec{B}) = \sum_{i=1}^n \hat{\varepsilon}_i^2 \rightarrow \min$

$\sum_{i=1}^n \hat{\varepsilon}_i^2 = \|\vec{\varepsilon}\|^2 = \|\vec{X} - Z^T \vec{B}\|^2$ - квадрат расстояния от точки \vec{X} до $Z^T \vec{B}$ в \mathbb{R}^n

$Z^T \vec{B}$ - точка линейного подпространства, порожденного векторами $Z^T \vec{t}$, где $\vec{t} \in \mathbb{R}^k$

Nota. Согласно Cond.1 размерность линейной оболочки, порожденной вектором $Z^T \vec{t}$, $\dim \langle Z^T \vec{t} \rangle = k$

Наименьшее расстояние получаем, когда квадрат расстояния от точки \vec{X} до данного подпространства, а вектор \vec{B} - проекция вектора \vec{X} на него

Таким образом, вектор $\vec{X} - Z^T \vec{B}$ должен быть ортогонален данному подпространству, то есть скалярное произведение вектора \vec{X} и всех векторов подпространства равно 0

$$(Z^T \vec{t}, \vec{X} - Z^T \vec{B}) = (Z^T \vec{t})^T (\vec{X} - Z^T \vec{B}) = \vec{t}^T (Z^T)^T (\vec{X} - Z^T \vec{B}) = \vec{t}^T Z (\vec{X} - Z^T \vec{B}) = \vec{t}^T (Z\vec{X} - ZZ^T \vec{B}) = 0 \quad \forall \vec{t} \in \mathbb{R}^k$$

Так как всем векторами подпространства ортогонален только нулевой вектор, то получаем, что $Z\vec{X} - ZZ^T \vec{B} = 0$ или $A\vec{B} = Z\vec{X}$ - нормальное уравнение (или система нормальных уравнений). Так как по свойству 2 матрица A невырожденная, то существует обратная, получаем решение системы: $\vec{B} = A^{-1}Z\vec{X}$

Свойства оценок метода наименьших квадратов

Добавим еще одно важное условие Cond.2: теоретические ошибки ε_i - независимы и имеют одинаковое нормальное распределение $N(0, \sigma^2)$. То есть $E\vec{\varepsilon} = \vec{0}$, $D\vec{\varepsilon} = \sigma^2 E_n$ (ковариации равны нулю в силу независимости)

Свойства:

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

$$\vec{B} - \vec{\beta} = A^{-1}Z\vec{X} - \vec{\beta} = A^{-1}Z(Z^T \vec{\beta} + \vec{\varepsilon}) - \vec{\beta} = A^{-1}Z\vec{\varepsilon}$$

$$2. \vec{B} - \text{несмещенная оценка для вектора } \vec{\beta}$$

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

$$3. \text{Матрица ковариаций } D\vec{B} = \sigma^2 A^{-1}$$

$$D\vec{B} = D(\vec{B} - \vec{\beta}) = D(A^{-1}Z\vec{\varepsilon}) = A^{-1}ZD\vec{\varepsilon}A^{-1T}Z^T = A^{-1}Z\sigma^2 E_n A^{-1T}Z^T = \sigma^2 A^{-1}(ZZ^T)A^{-1} = \sigma^2 A^{-1}$$

Следствие: дисперсии оценок b_i можно выразить через σ^2 и коэффициенты матрицы A^{-1} : $Db_i = \sigma^2 (A^{-1})_{ii}$

Оценка дисперсии случайного члена

Обозначим $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2$. Ясно, что $\hat{\sigma}^2$ - точечная оценка неизвестной дисперсии σ^2 , однако она является смещенной оценкой

Пусть выполнены Cond.1 и Cond.2, тогда $\frac{n\hat{\sigma}^2}{\sigma^2} \in H_{n-k}$ и не зависит от \vec{B}

Так как $\frac{n\hat{\sigma}^2}{\sigma^2} \in H_{n-k}$, то $E\hat{\sigma}^2 = \frac{\sigma^2}{n} E \frac{n\hat{\sigma}^2}{\sigma^2} = \frac{\sigma^2}{n} (n-k) = \frac{n-k}{n} \sigma^2 < \sigma^2$ - смещенная вниз оценка

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

Лекция 12.

Построение и анализ уравнения множественной линейной регрессии

Постановка задачи: пусть выявлена зависимость результата X от фактора Z_1, Z_2, \dots, Z_k

При $n > k$ получены результаты экспериментов $\vec{X} = (X_1, \dots, X_n)$ при соответствующих значениях факторов $\vec{Z}^{(j)} = (Z_1^{(j)}, Z_2^{(j)}, \dots, Z_k^{(j)})$

Предполагаем, что зависимость всех факторов от X линейная. Требуется по данным построить линейную модель, наилучшим образом объясняющую предсказывающую поведение X

Мультиколлинеарность

Def. Мультиколлинеарность - наличие линейной связи между всеми или несколькими факторами. Неприятные последствия:

- Оценки параметров становятся ненадежными, имеют большие стандартные ошибки и малую значимость. Небольшое изменение данных приводит к заметному изменению оценок
- Трудно определить изолированное влияние конкретного фактора на результат и выявить смысл данного влияния

Ex. Исследовалась зависимость веса от роста и размера обуви. В одной группе студентов модель оказалась такая: $X - \bar{x} = 0.9(Z_1 - \bar{Z}_1) + 0.1(Z_2 - \bar{Z}_2)$

В первой группе уравнение получилось таким: $X - \bar{x} = 0.9(Z_1 - \bar{Z}_1) + 0.1(Z_2 - \bar{Z}_2)$

А во второй группе таким: $X - \bar{x} = 0.2(Z_1 - \bar{Z}_1) + 0.8(Z_2 - \bar{Z}_2)$

Отбор факторов в уравнении регрессии

Чтобы избавиться от неприятных последствий, нужно отобрать факторы, которые мы в дальнейшем будем учитывать. Для этого строим корреляционную матрицу, состоящую из коэффициентов корреляции между результатом, факторами и факторов между собой

$$r = \begin{pmatrix} 1 & r_{X,Z_1} & r_{X,Z_2} & \dots & r_{X,Z_k} \\ r_{X,Z_1} & 1 & r_{Z_1,Z_2} & \dots & r_{Z_1,Z_k} \\ r_{X,Z_2} & r_{Z_2,Z_1} & 1 & \dots & r_{X,Z_k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_{X,Z_k} & r_{Z_k,Z_1} & r_{Z_k,Z_2} & \dots & 1 \end{pmatrix}$$

Алгоритм следующий:

1. Первым берем фактор, имеющий наибольшую корреляцию с результатом
2. Затем добавляем факторы, которые с одной стороны имеют наибольшую корреляцию с результатом, а с другой стороны наименьшую корреляцию с уже имеющимися факторами

Ех. Для данной корреляционной матрицы лучше всего брать Z_2 (фактор имеет лучшую корреляцию с X), а затем Z_3 (этот фактор меньше всего коррелирует с Z_2)

	X	Z ₁	Z ₂	Z ₃
X	1	-	-	-
Z ₁	0.81	1	-	-
Z ₂	0.85	0.93	1	-
Z ₃	-0.65	-0.38	-0.28	1

Анализ уравнения линейной регрессии

Пусть $X = \beta_0 + \beta_1 Z_1 + \dots + \beta_k Z_k + \varepsilon$ - теоретическая модель линейной регрессии, $\varepsilon \in N(0, \sigma^2)$ и независимы

По методу наименьших квадратов была построена модель $\hat{X} = b_0 + b_1 Z_1 + \dots + b_k Z_k$. Тогда $\hat{\varepsilon}_i = X_i - \hat{X}_i$ - экспериментальная (или эмпирическая) ошибка

Согласно следствию из теоремы $S^2 = \frac{1}{n-k-1} \sum_{i=1}^n \hat{\varepsilon}_i^2$ - несмещенная оценка σ^2

Def. Величину S называют стандартной ошибкой регрессии

Из прошлых лекций знаем, что $Db_i = \sigma^2(A^{-1})_{ii}$, где $A = ZZ^T$, $Z = (Z : (i)_j)$ - матрица плана

Из этого $\hat{D}b_i = S^2(A^{-1})_{ii}$ - оценка дисперсии

Def. $S_{b_i} = S\sqrt{(A^{-1})_{ii}}$ - стандартная ошибка коэффициента регрессии b_i

Уравнение регрессии стандартных масштабов

Nota. При обычном уравнении линейной регрессии трудно сравнить влияние различных факторов, так как они имеют разную природу и разные единицы измерения. Поэтому делают стандартизацию данных.

Пусть есть выборка $\vec{X} = (X_1, \dots, X_n)$ объема n . Тогда стандартизованный вектор \vec{t}_x состоит из значений $\frac{X_i - \bar{x}}{\hat{\sigma}_x}$

Получаем новые данные, которые можем считать новой выборкой стандартизированной случайной величины $\frac{X - EX}{\sigma_x}$ и которая не имеет единиц измерения

Свойства стандартизированных данных:

1. $\overline{t_x} = 0$
2. $D_{t_x}^* = 1$
3. $r_{x,y} = r_{t_x,t_y} = \overline{t_x t_y}$

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

$$r_{t_x,t_y} = \overline{t_x t_y}, \text{ так как } \overline{t_x} = \overline{t_y} = 0, D^* t_x = D^* t_y = 1$$

В уравнении регрессии данные результата и факторов заменяем на стандартизированные:

$$t_j = \left(\frac{Z_i^{(j)} - \overline{Z^{(j)}}}{\sigma_Z^{(j)}} \right) - \text{стандартизация } j\text{-ого фактора}, t_x = \left(\frac{X_i - \bar{X}}{\sigma_X} \right) - \text{стандартизация результата}$$

Так как стандартизация - линейная операция, то при соответствующей замене в уравнении получаем линейное уравнение, называемое уравнением регрессии стандартных масштабов:

$$t_x = \gamma_1 t_1 + \gamma_2 t_2 + \dots + \gamma_k t_k$$

Nota. Заметим, что γ_0 (свободный член) равен 0: $\overline{t_j} = 0$, поэтому линия уравнения пройдет через начало координат. При этом система нормальных уравнений приобретает более простой и наглядный вид:

$$\begin{cases} \gamma_1 + r_{Z_1,Z_2} \gamma_2 + \dots + r_{Z_1,Z_k} \gamma_k = r_{Z_1,x} \\ r_{Z_1,Z_2} \gamma_1 + \gamma_2 + \dots + r_{Z_2,Z_k} \gamma_k = r_{Z_2,x} \\ \dots \\ r_{Z_k,Z_1} \gamma_1 + r_{Z_k,Z_2} \gamma_2 + \dots + \gamma_k = r_{Z_k,x} \end{cases}$$

Или в матричной форме: $r\Gamma = r_x$, где r - корреляционная матрица, Γ - столбец коэффициентов регрессии, r_x - столбец из коэффициентов корреляции факторов с результатом

Нормальное уравнение регрессии: $A\vec{B} = Z\vec{X}$, где $A = ZZ^T$, Z - матрица плана

Допустим, что все данные стандартизированы, тогда для i -ого элемента

$$\begin{pmatrix} Z_1^{(i)} & \dots & Z_n^{(i)} \end{pmatrix} \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix} = \sum_{j=1}^n Z_j^{(i)} X_j = n \overline{Z^{(i)} X} = n r_{Z_i, X}$$

$$a_{ij} = \begin{pmatrix} Z_1^{(i)} & \dots & Z_n^{(i)} \end{pmatrix} \begin{pmatrix} Z_1^{(j)} \\ \vdots \\ Z_n^{(j)} \end{pmatrix} = \begin{cases} n \overline{Z_i Z_j} = \gamma_1 r_{Z_i, Z_j}, & \text{при } i \neq j \\ n \overline{Z_i^2} = n D Z_i = n, & \text{при } i = j \end{cases}$$

Перевод коэффициентов можно выполнять по формулам:

$$b_i = \gamma_i \frac{\sigma_X}{\sigma_{Z_i}}; \quad \gamma_i = b \frac{\sigma_{Z_i}}{\sigma_X}, \quad b_0 = \bar{x} - \sum_{i=1}^n b_i \bar{z}_i$$

Ex. Уравнение парной линейной регрессии: $\frac{x - \bar{x}}{\hat{\sigma}_x} = \hat{r} \frac{z - \bar{z}}{\hat{\sigma}_z}$

В стандартных масштабах получаем $t_x = \hat{r}t_z$

Тогда $y_1 = \hat{r}$

Смысл стандартизированных коэффициентов

Значение y_i можно трактовать как величину прямого влияния i -ого фактора на результат. Остальные слагаемые в i -ом уравнении можно трактовать как величины косвенного влияния остальных факторов

$$r_{Z_1 Z_1} y_1 + \dots + y_i + \dots + r_{Z_i Z_k} y_k = r_{Z_i X}$$

Величину $r_{Z_i X}$ можно рассматривать как сумму прямого и косвенного влияний

Nota. Для измерения тесноты отдельной связи между отдельным фактором и результатом при очищении влияния других факторов есть понятие коэффициента частной корреляции.

При $k = 2$:
$$r_{X, Z_1 / Z_2} = \frac{r_{X, Z_1} - r_{X, Z_2} r_{Z_1, Z_2}}{\sqrt{(1 - r_{X, Z_2}^2)(1 - r_{Z_1, Z_2}^2)}}$$

Коэффициенты детерминации и множественной корреляции

Допустим, что как и в случае парной линейной регрессии дисперсию результата X можно разложить на объясненную и необъясненную $D(X) = D(\hat{X}) + D(\hat{\epsilon})$

Nota. В зарубежной литературе используется $RSS = D(\hat{\epsilon}) = \sum \hat{\epsilon}_i^2$

Def. Коэффициентом детерминации R^2 называется величина $R^2 = 1 - \frac{D(\hat{\epsilon})}{D(X)}$

Его можно трактовать как долю объясненной дисперсии, а $1 - R^2$ - как долю необъясненной

Свойства:

1. $0 \leq R^2 \leq 1$
2. Если $R^2 = 1$, то $D(\hat{\epsilon}) = 0$, то есть все данные лежат в гиперплоскости построенного уравнения регрессии
3. Если $R^2 = 0$, то $D(\hat{X}) = 0$, то есть все $\hat{X}_i = \bar{x} \implies b_1 = b_2 = \dots = b_k = 0$, то есть модель ничего не объясняет
4. Чем больше R^2 , тем лучше
5. В случае линейного уравнения регрессии $R^2 = \sum_{i=1}^k y_i r_{Z_i, X}$

Def. Величина R называется коэффициентом множественной корреляции, который показывает силу линейной связи

Скорректированный коэффициент детерминации

При добавлении в модель новых факторов R^2 как правило увеличивается, хотя не всегда есть смысл добавлять новые факторы. Для выяснения того, следует ли это делать, используется скорректированный коэффициент детерминации $\overline{R^2} = 1 - \frac{k-1}{n-k-1} \frac{D(\hat{\varepsilon})}{D(X)}$, где n - число экспериментов, а k - число факторов

Проверка гипотез по значимости уравнения регрессии

- а) F-тест: проверка гипотезы о значимости уравнения регрессии в целом $H_0 : R_{\text{теор}}^2 = 0$ (уравнение регрессии статистически незначимо) против $H_1 : R_{\text{теор}}^2 \neq 0$

Th. Если $H_0 : R_{\text{теор}}^2 = 0$ верна, то $F = \frac{R^2}{1-R^2} \frac{n-k-1}{k} \in F(k, n-k-1)$

Получаем критерий, называемый F-тестом: t_α - квантиль $F(k, n-k-1)$ уровня значимости

$$\begin{cases} H_0 : R_{\text{теор}}^2 = 0, & \text{если } F < t_\alpha \\ H_0 : R_{\text{теор}}^2 \neq 0, & \text{если } F \geq t_\alpha \end{cases}$$

- б) T-тест: проверка гипотезы о значимости отдельного коэффициента регрессии $H_0 : \beta_i = 0$ (b_i статистически незначим) против $H_1 : \beta_i \neq 0$

Th. Если $H_0 : \beta_i = 0$ верна, то $T_i = \frac{b_i}{S_{b_i}} \in T_{n-k-1}$

Получаем критерий, называемый T-тестом: t_α - квантиль двухстороннего распределения Стьюдента $|T_{n-k-1}|$ уровня значимости α

$$\begin{cases} H_0 : \beta_i = 0, & \text{если } |T_i| < t_\alpha \\ H_0 : \beta_i \neq 0, & \text{если } |T_i| \geq t_\alpha \end{cases}$$

Nota. T-тест служит для отсева несущественных факторов из модели при условии, что все другие факторы включены в модель

Nota. При мультиколлинеарности возможно, что уравнение имеет высокую значимость, а большинство коэффициентов не проходит T-тест

Nota. При применении T-теста убираем только один фактор, далее строим новую модель и для нее опять проводим T-тест. Удаление 2 факторов может привести к неопределенным результатам

Лекция 13.

Невазия регрессионного анализа

Пусть имеется уравнение общей линейной регрессии $\vec{X} = Z^T \vec{\beta} + \vec{\varepsilon}$, где n - число экспериментов, \vec{X} - столбец результатов экспериментов, Z - матрица плана, $\vec{\beta}$ - столбец коэффициентов регрессии, $\vec{\varepsilon}$ - вектор теоретических ошибок

При этом ранее предполагали, что выполнены условия:

1. Cond. 1: Строки Z - независимы
2. Cond. 2: $\varepsilon_i \in N(0, \sigma^2)$ и независимы

Условие 2 часто нарушается

Взвешенный метод наименьших квадратов

Пусть $\varepsilon_i \in N(0, v_i \sigma^2)$ и независимы (то есть дисперсия ошибки зависит от номера наблюдения). Другими словами, $D\vec{\varepsilon} = \sigma^2 V$, где $V = \text{diag}(v_1, \dots, v_n)$

Логично наблюдениям с меньшей дисперсии ошибки предать больший вес. Пусть вес $w_i = \frac{1}{v_i}$.

Домножим обе части уравнения регрессии на $\sqrt{w_i}$, тогда получим $\tilde{X} = \tilde{Z}^T \vec{\beta} + \tilde{\varepsilon}$, где $\tilde{x}_i = \sqrt{w_i} x_i$, $\tilde{Z}_i^{(j)} = \sqrt{w_i} Z_i^{(j)}$, $\tilde{\varepsilon}_i = \sqrt{w_i} \varepsilon_i$

$D\vec{\varepsilon} = D(\sqrt{w_i} \varepsilon_i) = w_i D\varepsilon_i = \frac{1}{v_i} v_i \sigma^2 = \sigma^2$ - получаем, что $D\vec{\varepsilon} = \sigma^2 E_n$, то есть стандартную ситуацию

Тогда оценки \vec{b} будут несмещенными и эффективными

Недостаток у этого метода: нужно знать коэффициенты v_i

Ex. Рассмотрим модель линейной парной регрессии без свободного члена $X = \beta_0 Z + \varepsilon$

Теоретическое уравнение $A\vec{B} = Z\vec{X}$, где $Z = (Z_1, \dots, Z_n)$ - матрица плана, $\vec{B} = \hat{\beta}_0$, $A = ZZ^T = z_1^2 + \dots + z_n^2$, $Z\vec{X} = z_1 x_1 + \dots + z_n x_n$

$$\sum_{i=1}^n z_i^2 \hat{\beta}_0 = \sum_{i=1}^n z_i x_i \implies \hat{\beta}_0 = \frac{\sum z_i x_i}{\sum z_i^2} - \text{оценка МНК}$$

По взвешенному методу наименьших квадратов $\tilde{\beta}_0 = \frac{\sum w_i z_i x_i}{\sum w_i z_i^2}$ - оценка взвешенного МНК

Ex. a. Взвешенное среднее

Допустим, что проводим серию измерений «скоропортящимся» инструментом. При $Z \equiv 1$: $X = \beta_0 + \varepsilon$, $\varepsilon_i \in N(0, v_i \sigma^2)$, $w_i = \frac{1}{v_i}$

Тогда $\hat{\beta}_0 = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$ - взвешенное среднее

Ex. б. Повторное наблюдение

Пусть было n серий по k_i наблюдений ($1 \leq i \leq n$). В каждой серии вычислили выборочное среднее \bar{x}_i . Если $\varepsilon \in N(0, \sigma^2)$, то дисперсия ошибки для каждого выборочного среднего $D\varepsilon_i = \frac{\sigma^2}{k_i}$.

Оценка по результатам всех наблюдений будет $\hat{\beta}_0 = \frac{\sum_{i=1}^n k_i \bar{x}_i}{\sum_{i=1}^n k_i}$

Ех. в. Пропорция

Пусть X - потери тепла в квартире. Основной фактор Z - разница температур снаружи и внутри. Так как при $Z = 0$ $X = 0$, то уравнением регрессии будет $X = \beta_0 Z + \varepsilon$

Логично предположить, что дисперсия ошибки зависит от величины Z . Рассмотрим две гипотезы:

1. Дисперсия ошибки прямо пропорциональна Z : $D\varepsilon_i = CZ_i = \sigma^2 \frac{CZ_i}{\sigma^2}$

$$\text{Тогда } w_i = \frac{\sigma^2}{CZ_i} \text{ и } \hat{\beta}_0 = \frac{\sum_{i=1}^n \frac{\sigma^2}{CZ_i} Z_i X_i}{\sum_{i=1}^n \frac{\sigma^2}{CZ_i} Z_i^2} = \frac{\sum X_i}{\sum Z_i} = \frac{\bar{x}}{\bar{z}}$$

2. Дисперсия ошибки квадратично зависит от Z : $D\varepsilon_i = CZ_i^2 = \sigma^2 \frac{CZ_i^2}{\sigma^2}$

$$\text{Тогда } w_i = \frac{\sigma^2}{CZ_i^2} \text{ и } \hat{\beta}_0 = \frac{\sum_{i=1}^n \frac{\sigma^2}{CZ_i^2} Z_i X_i}{\sum_{i=1}^n \frac{\sigma^2}{CZ_i^2} Z_i^2} = \frac{\sum X_i}{\sum Z_i} = \frac{\sum \frac{X_i}{Z_i}}{n} = \overline{\left(\frac{x}{z}\right)}$$

Коррелированные наблюдения

Пусть ошибки не только имеют различные дисперсии, но и коррелированы между собой: $\text{cov}(\varepsilon_i, \varepsilon_j) = v_{ij}$

Тогда $D\vec{\varepsilon} = \sigma^2 V$, где $V = (v_{ij})$

Так как матрица ковариаций симметричная и положительно определенная, то существует \sqrt{V} .

Домножим обе части уравнения регрессии на $\sqrt{V^{-1}}$:

$$\vec{X} = Z^T \vec{\beta} + \vec{\varepsilon} \quad \Big| \cdot \sqrt{V^{-1}}$$

$$\vec{X} = \tilde{Z}^T \vec{\beta} + \vec{\varepsilon}, \text{ где } \vec{X} = \sqrt{V^{-1}} \vec{X}, \tilde{Z} = \sqrt{V^{-1}} Z, \vec{\varepsilon} = \sqrt{V^{-1}} \vec{\varepsilon}$$

Тогда матрица ковариаций нового вектора ошибок будет $D\vec{\varepsilon} = D(\sqrt{V^{-1}} \vec{\varepsilon}) = \sqrt{V^{-1}} D\vec{\varepsilon} (\sqrt{V^{-1}})^T = \sigma^2 I_n$

То есть получили классическую ситуацию, когда выполнено Cond. 2 и вектор оценок $\hat{\beta}_0$ будет несмещенным и эффективным

Составление матрицы плана при управляемом эксперименте

Если строки матрицы плана взять ортогональными, то дисперсии оценки коэффициентов b_i регрессии будут минимальными. Поэтому лучше матрицу плана составлять таким образом:

Дисперсии оценок при этом $Db_i = \sigma^2 A_{ii}^{-1}$. Если Z - ортогональная (не обязательно нормированная), то $A = ZZ^T = E_n$, а $Db_i = \frac{\sigma^2}{Z_i^2}$. Несложно доказать, что во всех других случаях дисперсия будет больше

Метод главных осей

Помимо метода наименьших квадратов существует метод «главных осей». Идея следующая: матрицу ковариаций приводим к диагональной форме

В МНК мы минимизируем расстояние отрезков, параллельных оси Оу, а в методе главных осей - перпендикуляр от точки до возможной прямой

Результатом метода главных осей получаем прямую, являющуюся главной осью эллипса, появляющегося из корреляционного облака

Нелинейные регрессии

Помимо общего МНК многие нелинейные зависимости могут быть сведены к линейным при помощи простых приемов

Ex. а. $X = \alpha + \beta f(Z) + \varepsilon$, где $f(x)$ - известная функция

Тогда можно взять новый фактор $Z' = f(Z)$, свели задачу к стандартной, получаем уравнение $X = \alpha + \beta Z' + \varepsilon$

Пример: $X = \alpha + \beta \ln Z + \varepsilon$, то $Z' = \ln Z$

Ex. б. $X = \alpha Z^\beta + \varepsilon$

Логарифмируем: $\ln X = \ln \alpha + \beta \ln Z + \ln \varepsilon \iff X' = \alpha' + \beta Z' + \varepsilon'$

Ex. в. $X = \alpha e^{\beta Z} + \varepsilon$

Логарифмируем: $\ln X = \ln \alpha + \beta Z + \ln \varepsilon \iff X' = \alpha' + \beta Z + \varepsilon'$

Ex. г. Зависимость в виде полинома: $X = \beta_0 + \beta_1 Z + \beta_2 Z^2 + \dots + \beta_k Z^k$

Введем новые факторы $Z_1 = Z, Z_2 = Z^2, \dots, Z_k = Z^k$

$X = \beta_0 + \beta_1 Z + \beta_2 Z_2 + \dots + \beta_k Z_k$

При этом, чтобы избежать мультиколлинеарность, лучше брать $k < 4$. При больших k получить многочлен большой степени, который сможет гарантировано пройти через все точки - это будет статистически незначимо

Nota. Если из теории мы знаем вид зависимости и подбираем ее под данные, то желательно строить модель как можно проще

Nota. Из построенных моделей предпочтительней та, где коэффициент детерминации больше

Построение даже удачной регрессионной модели не означает появление причинно-следственной связи. Исторический пример: исследовалась точность бомбометания от различных условий.

Пусть X - точность, Z_1 - высота, Z_2 - ветер, Z_3 - количество истребителей противника

Построили модель $X = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2 + \beta_3 Z_3$ и получили $\hat{\beta}_3 < 0$ - то есть при большем числе техники противника точность увеличивалась. Оказалось, что не был учтен фактор облачности - при нем $\hat{\beta}_3 > 0$ и коэффициент детерминации улучшился

Или другой пример: корреляция численности аистов и рождения детей в Голландии XX века оказалась прямой

Х. Программа экзамена в 2024/2025

1. Выборочная функция распределения. Теоремы Гливенко-Кантелли и Колмогорова.
2. Вариационный и интервальный вариационный ряды. Полигон и гистограмма, ее свойства.
3. Точечные оценки. Их свойства: состоятельность, несмещенность, эффективность.
4. Точечные оценки моментов. Свойства оценок математического ожидания и дисперсии.
5. Метод моментов. Пример.
6. Метод максимального правдоподобия. Пример.
7. Информация Фишера. Неравенство Рао-Крамера (без док-ва).
8. Основные распределения математической статистики: хи-квадрат, Стьюдента, Фишера-Снедекора. Их свойства.
9. Линейные преобразования нормальных выборок. Теорема об ортогональном преобразовании.
10. Лемма Фишера.
11. Основная теорема о связи точечных оценок нормального распределения и основных распределений статистики.
12. Квантили распределений (оба определения). Функции для их вычисления в EXCEL.
13. Интервальные оценки. Определения, смысл, терминология.
14. Доверительный интервал для математического ожидания нормального распределения при известном σ .
15. Доверительный интервал для математического ожидания нормального распределения при неизвестном σ .
16. Доверительный интервал для дисперсии нормального распределения при неизвестном σ .
17. Доверительный интервал для дисперсии нормального распределения при известном σ .
18. Проверка статистических гипотез. Определения, терминология. Уровень значимости и мощность критерия.
19. Способы сравнения критериев проверки гипотез.
20. Построение критериев согласия (основные принципы).
21. Гипотеза о среднем нормальной совокупности с известной дисперсией.
22. Гипотеза о среднем нормальной совокупности с неизвестной дисперсией.
23. Доверительные интервалы как критерии гипотез о параметрах распределения.
24. Критерий хи-квадрат для параметрической гипотезы.
25. Критерий хи-квадрат для гипотезы о распределении.
26. Критерий Колмогорова для гипотезы о распределении.
27. Критерий Колмогорова-Смирнова.
28. Критерий Фишера.
29. Критерий Стьюдента.
30. Понятие статистической зависимости. Корреляционное облако и корреляционная таблица.

Первоначальные выводы.

31. Критерий хи-квадрат для проверки независимости.
32. Однофакторный дисперсионный анализ. Общая, межгрупповая и внутригрупповая дисперсии. Теорема о разложении дисперсии.
33. Однофакторный дисперсионный анализ. Проверка гипотезы о влиянии фактора.
34. Математическая модель регрессии. Основные понятия и определения. Метод наименьших квадратов.
35. Вывод уравнения линейной парной регрессии. Геометрический смысл прямой регрессии.
36. Выборочный коэффициент линейной корреляции. Проверка гипотезы о его значимости.
37. Выборочное корреляционное отношение, его свойства.
38. Свойства ошибок в модели линейной парной регрессии. Анализ дисперсии фактора-результата. Коэффициент детерминации, его свойства.
39. Проверка гипотезы о значимости уравнения линейной регрессии. Связь между коэффициентом детерминации и коэффициентом линейной корреляции.
40. Теорема Гаусса-Маркова.
41. Стандартные ошибки коэффициентов регрессии. Их доверительные интервалы.
42. Прогнозирование в модели линейной парной регрессии. Стандартная ошибка прогноза, доверительный интервал прогноза.
43. Общая модель линейной регрессии. Вывод нормального уравнения.
44. Свойства ОНМК в уравнении общей линейной регрессии.
45. Основная теорема об ОМНК (п.2 без доказательства).
46. Мультиколлинеарность, ее неприятные последствия. Основные принципы отбора факторов в модель общей линейной регрессии.
47. Стандартная ошибка общей линейной регрессии и стандартные ошибки коэффициентов регрессии. Проверка гипотезы о значимости отдельного коэффициента регрессии.
48. Уравнение регрессии в стандартных масштабах. Смысл стандартизованных коэффициентов. Разложение влияния фактора на прямое и косвенное.
49. Коэффициенты детерминации и множественной корреляции, их свойства. Проверка гипотезы о значимости уравнения регрессии в целом.
50. Взвешенный МНК.
51. Приемы сведения нелинейных регрессий к линейным.
52. Математические датчики случайных чисел.
53. Моделирование случайных величин методом обратной функции (включая дискретный случай).
54. Моделирование нормальной случайной величины.
55. Быстрый показательный датчик.
56. Моделирование дискретных случайных величин.
57. Метод Монте-Карло. Общая постановка, оценка погрешности.

58. Вычисление определенного и кратного интегралов методом Монте-Карло. Метод расслоенной выборки.