

Статистика

Конспект курса

Мат-Мех, ПМИ

5–7 семестры (2015–2016)

\$Revision: 1.52 \$

Содержание

1. Описательная статистика	3
1.1. Выборка и эмпирическая случайная величина	3
1.2. Виды признаков	4
2. Точечное оценивание	4
2.1. Характеристики распределений и метод подстановки	4
2.2. Моменты	5
2.2.1. Характеристики положения	5
2.2.2. Характеристики разброса	5
2.2.3. Анализ характера разброса	7
2.3. Свойства оценок	7
2.3.1. Несмещенность	7
2.3.2. Состоятельность	8
2.3.3. Асимптотическая нормальность	10
2.3.4. Эффективность	10
2.4. Метод моментов	10
2.5. Метод оценки максимального правдоподобия	11
2.5.1. Функция правдоподобия и метод	11
2.5.2. Информационное количество Фишера и эффективность оценки MLE	12
3. Некоторые распределения, связанные с нормальным	13
3.1. Распределение $\chi^2(m)$	13
3.2. Распределение Стьюдента $t(m)$	14
3.3. Распределение Фишера	14
3.4. Квадратичные формы от нормально распределенных случайных величин	14
4. Проверка гипотез	15
4.1. Построение критерия	15
4.1.1. Понятие гипотезы и критерия	15
4.1.2. Построение критерия при помощи статистики критерия	17
4.1.3. Схема построения критерия с помощью статистики	19
4.2. Проверка гипотезы о значении мат. ожидания (t -критерий)	21
4.2.1. $D\xi = \sigma^2 < \infty$	21
4.2.2. $D\xi$ неизвестна	22
4.3. Проверка гипотезы о значении дисперсии в нормальной модели (критерий χ^2)	23
4.3.1. $E\xi = a < \infty$	23

4.3.2. $E\xi$ неизвестно	23
4.4. Критерий χ^2 согласия с видом распределения	24
4.4.1. Распределение с известными параметрами	24
4.4.2. Распределение с неизвестными параметрами	25
4.4.3. Согласие с нормальным распределением	26
4.5. Критерий Колмогорова-Смирнова согласия с видом распределения	27
4.5.1. Произвольное абсолютно непрерывное распределение	27
4.5.2. Нормальное распределение	28
4.6. Критерий типа ω^2	28
4.7. Визуальное определение согласия с распределением	29
4.7.1. P-P plot	29
4.7.2. Q-Q plot	29
4.8. Гипотеза о равенстве распределений	29
4.9. Равенство математических ожиданий для независимых выборок	30
4.9.1. Двухвыборочный t -критерий	30
4.9.2. Непараметрический t -критерий	32
4.9.3. Критерии суммы рангов Wilcoxon	32
4.9.4. Критерий Mann-Whitney (U test)	32
4.9.5. Критерий серий (runs)	33
4.9.6. Критерий равенства распределений	33
4.10. Равенство математических ожиданий для парных (зависимых) выборок	33
4.10.1. t -критерий	33
4.10.2. Непараметрический тест знаков (Sign test)	34
4.10.3. Непараметрический критерий (Paired Wilcoxon; Wilcoxon signed-rank test)	34
4.11. Равенство дисперсии для двух распределений	34
4.11.1. Критерий Фишера	34
4.11.2. Критерий Левена (Levene's test)	35
4.11.3. Критерий Brown-Forsythe	35
5. Доверительное оценивание	35
5.1. Мотивация и определение	35
5.2. Доверительные интервалы для математического ожидания и дисперсии в нормальной модели	35
5.2.1. Доверительный интервал для a	35
5.2.2. Доверительный интервал для σ^2	36
5.3. Асимптотический доверительный интервал для математического ожидания в модели с конечной дисперсией	37
5.4. Асимптотический доверительный интервал для параметра на основе MLE	37
5.5. Доверительный интервал для проверки гипотезы о значении параметра	38
5.5.1. Использование SE для построения доверительных интервалов	38
6. Корреляционный и регрессионный анализы	38
6.1. Вероятностная независимость	38
6.1.1. Визуальное определение независимости	38
6.1.2. Критерий независимости χ^2	39
6.2. Линейная / полиномиальная зависимость	40
6.3. Метод наименьших квадратов (Ordinary Least Squares)	41
6.4. Корреляционное отношение	41
6.5. Регрессия	42
6.5.1. Парная линейная регрессия	42
6.5.2. Модель линейной регрессии	43
6.5.3. Доверительные интервалы для β_1 и β_2	44
6.6. Частная корреляция	45

6.7. Зависимость между порядковыми признаками	46
6.7.1. Ранговый коэффициент Спирмана	46
6.7.2. Ранговый коэффициент Кэндалла $\tau(\xi, \eta)$	48
6.8. Корреляционные матрицы	49
А. Другие полезные распределения случайных величин	50
А.1. Пуассона	50
А.2. Логнормальное	50
В. Свойства условного математического ожидания	50

1. Описательная статистика

1.1. Выборка и эмпирическая случайная величина

Пусть $\xi \sim \mathcal{P}$ — случайная величина с распределением \mathcal{P} .

Определение. Повторной независимой выборкой объема n (до эксперимента) называется набор

$$\mathbf{x} = (x_1, \dots, x_n), \quad x_i \sim \mathcal{P} \quad \forall i \in 1:n, \quad x_1 \perp \dots \perp x_n$$

независимых в совокупности одинаково распределенных случайных величин с распределением \mathcal{P} .

Определение. Повторной независимой выборкой объема n (после эксперимента) называется набор реализаций, т.е. конкретных значений ξ , случайных величин x_i :

$$\mathbf{x} = (x_1, \dots, x_n), \quad x_i \in \text{im } \xi \quad \forall i \in 1:n.$$

Определение. Эмпирической случайной величиной $\hat{\xi}_n$ называется случайная величина с дискретным распределением

$$\hat{\xi}_n \sim \hat{\mathcal{P}}_n : \begin{pmatrix} x_1 & \dots & x_n \\ 1/n & \dots & 1/n \end{pmatrix}.$$

Замечание. Подходящее определение выбирается по контексту.

Если ξ имеет дискретное распределение, то выборку можно *сгруппировать*; тогда получим случайную величину $\hat{\xi}_m$ с распределением

$$\hat{\mathcal{P}}_m : \begin{pmatrix} x_1^* & \dots & x_m^* \\ \omega_1 & \dots & \omega_m \end{pmatrix} \quad \omega_i = \frac{\nu_i}{n},$$

где x_i^* — уникальные значения из выборки \mathbf{x} , а ν_i — число x_i^* в \mathbf{x} (т.н. «абсолютная частота»; тогда ω_i — «относительная частота»). В противном случае, можно разбить интервал всевозможных значений выборки на m подынтервалов: $\{[e_0, e_1), \dots, [e_{m-1}, e_m)\}$ и считать число наблюдений $\nu_i = \nu_i[e_{i-1}, e_i)$, попавших в интервал.

Следствие. По ЗБЧ (теореме Бернулли),

$$\omega_i \xrightarrow{P} p_i = P(e_{i-1} \leq \xi < e_i),$$

т.е. относительная частота является хорошей оценкой вероятности на больших объемах выборки.

1.2. Виды признаков

Виды признаков случайной величины $\xi : (\Omega, \mathcal{F}, P) \rightarrow (V, \mathfrak{A})$ характеризуются тем, что из себя представляет множество V и что можно делать с его элементами.

Количественные признаки: $V \subset \mathbb{R}$

По типу операций:

- Аддитивные: заданы, т.е. имеют смысл в контексте данного признака, операции $+$, $-$
- Мультипликативные: заданы операции \cdot , $/$; признак принимает не отрицательные значения.

По типу данных:

- Непрерывные
- Дискретные

Порядковые признаки V — упорядоченное множество, определены отношения $>$, $=$.

Качественные признаки на V заданы отношения $=$, \neq

Пример. Цвет глаз, имена, пол.

2. Точечное оценивание

2.1. Характеристики распределений и метод подстановки

Определение. *Статистика* — измеримая функция от выборки.

Обобщением статистики является понятие характеристики.

Определение. *Характеристика* — функционал от распределения:

$$T : \{\mathcal{P}\} \rightarrow V.$$

Где V — измеримое пространство, чтобы на нём можно было завести σ -алгебру.

Замечание. Чаще всего, $V = \mathbb{R}$.

Определение. Выделяют *генеральные* характеристики $T(\mathcal{P}) =: \theta$ и *выборочные* характеристики $T(\hat{\mathcal{P}}_n)$. *Оценка* — выборочная характеристика $T(\hat{\mathcal{P}}_n) =: \hat{\theta}_n$, не зависящая от генеральной характеристики θ .

Следствие. *Выражения для вычисления генеральных и выборочных характеристик отличаются только используемыми мерами (\mathcal{P} и $\hat{\mathcal{P}}_n$ соответственно).*

Определение. Пусть $\hat{\mathcal{P}}_n$ — распределение эмпирической случайной величины. Тогда *эмпирическая функция распределения* есть

$$\widehat{\text{cdf}}_{\xi}(x) = \text{cdf}_{\xi_n}(x) = \hat{\mathcal{P}}_n((-\infty, x)) = \int_{-\infty}^x d\hat{\mathcal{P}}_n = \sum_{x_i: x_i \leq x} \frac{1}{n} = \frac{|\{x_i \in \mathbf{x} : x_i \leq x\}|}{n}.$$

Утверждение. Пусть $\widehat{\text{cdf}}_{\xi}$ — эмпирическая функция распределения, cdf_{ξ} — функция распределения ξ . Тогда, по теореме Гливенко-Кантелли,

$$\sup_x \left| \widehat{\text{cdf}}_{\xi}(x) - \text{cdf}_{\xi}(x) \right| \xrightarrow{\text{a.s.}} 0.$$

Замечание. Поскольку $\widehat{\text{cdf}}_{\xi}(x) = \omega_x$, где ω_x — частота попадания наблюдений в интервал в $(-\infty, x)$, а $\text{cdf}_{\xi}(x) = P(\xi \in (-\infty, x))$ — вероятность того же события, то можно применить теорему Бернулли (ЗБЧ):

$$\widehat{\text{cdf}}_{\xi}(x) \xrightarrow{P} \text{cdf}_{\xi}(x).$$

Следствие. *Значит, при достаточно больших n , в качестве интересующей характеристики θ распределения \mathcal{P} можем брать ее оценку $\hat{\theta}_n$ — аналогичную характеристику $\hat{\mathcal{P}}_n$.*

2.2. Моменты

Определение. Генеральные и соответствующие им выборочные характеристики k -го момента и k -го центрального момента:

$$\begin{aligned} m_k &= \int_{\mathbb{R}} x^k d\mathcal{P} & \hat{m}_k &= \int_{\mathbb{R}} x^k d\hat{\mathcal{P}}_n = \frac{1}{n} \sum_{i=1}^n x_i^k \\ m_k^{(0)} &= \int_{\mathbb{R}} (x - m_1)^k d\mathcal{P} & \hat{m}_k^{(0)} &= \int_{\mathbb{R}} (x - \hat{m}_1)^k d\hat{\mathcal{P}}_n = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{m}_1)^k. \end{aligned}$$

2.2.1. Характеристики положения

В качестве характеристики положения выделяется 1-й момент — математическое ожидание и выборочное среднее:

$$m_1 = E\xi, \quad \hat{m}_1 =: \bar{x} = \widehat{E\xi} = E\hat{\xi}_n.$$

Замечание. В случае мультипликативных признаков можно посчитать среднее геометрическое; часто логарифмируют и считают среднее арифметическое.

Определение. Пусть $p \in [0, 1]$ и $\text{cdf} = \text{cdf}_{\mathcal{P}}$. p -квантилью называется

$$\text{qnt}_{\mathcal{P}}(p) =: z_p = \sup \{z : \text{cdf}(z) \leq p\}.$$

Квартиль есть 1/4- или 3/4-квантиль.

Определение. Медиана есть 1/2-квантиль:

$$\text{med } \xi = z_{1/2}.$$

Определение. Мода ($\text{mode } \xi$) есть точка локального максимума плотности.

По методу подстановки можем получить аналогичные выборочные характеристики.

Определение. Выборочная p -квантиль есть такая точка \hat{z}_p , что она больше по значению $|\mathbf{x}| \cdot p = np$ точек из выборки:

$$\hat{z}_p = \sup \left\{ z : \widehat{\text{cdf}}_{\xi}(z) \leq p \right\} = x_{[np]+1}.$$

Определение. Выборочная медиана упорядоченной выборки $\mathbf{x} = (x_{(1)}, \dots, x_{(n)})$ есть

$$\hat{z}_{1/2} = \widehat{\text{med}} = \begin{cases} x_{(k+1)} & n = 2k + 1 \\ \frac{x_{(k)} + x_{(k+1)}}{2} & n = 2k \end{cases}$$

Определение. Выборочная мода ($\widehat{\text{mode}}$) есть значение из выборки, которое чаще всего встречается.

2.2.2. Характеристики разброса

В качестве характеристики разброса выделяется 2-й центральный момент — дисперсия и выборочная дисперсия:

$$m_2^{(0)} = D\xi \quad \hat{m}_2^{(0)} =: s^2 = \widehat{D\xi} = D\hat{\xi}_n = \begin{cases} E(\hat{\xi}_n - E\hat{\xi}_n)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\ E\hat{\xi}_n^2 - (E\hat{\xi}_n)^2 = \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right) - \bar{x}^2. \end{cases}$$

Замечание. Если среднее $E\xi = a$ известно, то дополнительно вводится

$$s_a^2 := \begin{cases} \frac{1}{n} \sum_{i=1}^n (x_i - a)^2 \\ \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right) - a^2. \end{cases}$$

Пример (Оценка дисперсии оценки мат. ожидания). Пусть строится оценка мат. ожидания \bar{x} . Может интересовать точность построенной оценки. Вычислим дисперсию теоретически, после чего оценим точность по выборке:

$$D\bar{x} = D\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n^2} \sum_{i=1}^n Dx_i = \frac{1}{n^2} \sum_{i=1}^n D\xi = \frac{D\xi}{n},$$

откуда

$$\widehat{D\bar{x}} = \frac{s^2}{n}.$$

Пример (Дисперсия оценки дисперсии). См. по ссылке¹.

Определение (Энтропия). Количество информации, необходимое для выявления объекта из n -элементного множества вычисляется по *формуле Хартли*:

$$H = \log_2 n$$

(множество это следует итеративно разбивать пополам, откуда и оценка). Пусть теперь множество не равновероятно, т.е. задано дискретное распределение

$$\mathcal{P}_\xi : \begin{pmatrix} x_1 & \dots & x_n \\ p_1 & \dots & p_n \end{pmatrix}.$$

Тогда количество информации $H(\xi)$, которую нужно получить, чтобы узнать, какой исход эксперимента осуществлен, вычисляется по формуле Шеннона и называется *энтропией*:

$$H(\xi) = \sum_{i=1}^n p_i \log_2 \frac{1}{p_i}.$$

Замечание. В случае равномерного дискретного распределения, конечно, $H = H(\xi)$.

Определение. *Выборочное стандартное отклонение* есть

$$SD := \sqrt{\widehat{D\xi}} = s.$$

Это показатель разброса случайной величины; показатель того, насколько элементы выборки отличаются от выборочного среднего по значению.

Определение. *Стандартная ошибка* оценки есть

$$SE := \sqrt{\widehat{D\hat{\theta}}}.$$

Это показатель разброса оценки случайной величины.

Замечание. В частном случае $\theta = E\xi$, $\hat{\theta} = \bar{x}$ получаем *выборочную стандартную ошибку среднего*

$$SE = \sqrt{\frac{D\xi}{n}} = \frac{s}{\sqrt{n}}.$$

Это, в свою очередь, показатель того, насколько выборочное среднее отличается от истинного.

Пример (Пример с мостом и машинами). FIXME

¹<http://mathworld.wolfram.com/SampleVarianceDistribution.html>

2.2.3. Анализ характера разброса

Определение. Коэффициент асимметрии Пирсона («скошенности»²⁾

$$\gamma_3 = A\xi = \frac{m_3^{(0)}}{\sigma^3} = \frac{E\xi - \text{med } \xi}{\sigma}.$$

Замечание. Не зависит от линейных преобразований.

Определение. Коэффициент эксцесса («крутизны», «kurtosis»):

$$\gamma_4 = K\xi = \frac{m_4^{(0)}}{\sigma^4} - 3.$$

Замечание. Величина $m_4^{(0)}/\sigma^4 = 3$ соответствует стандартному нормальному распределению. Так что можно сравнивать выборку и $\gamma_4 N(0, 1)$.

Замечание. При замене $z := (\xi - E\xi)/\sigma$ величину $m_4^{(0)}/\sigma^4 = E(z^4)$ можно интерпретировать как ожидание четвертой степени центрированных и нормированных данных. Точки выборки, лежащие внутри $E\xi \pm \sigma$ из-за малости по модулю не будут увеличивать значение коэффициента, в то время как аутлаеры будут или «тяжелые хвосты» плотности распределения будут. Поэтому γ_4 принимает большие значения на распределениях с «тяжелыми хвостами» или выборках с некоторым количеством аутлаеров.

Замечание. Справедлива оценка

$$\gamma_3^2 + 4 \leq \gamma_4 + 3 \leq \infty,$$

где минимум достигается $\text{Ber}(1/2)$.

Определение. Пусть $(\xi_1, \xi_2) \sim \mathcal{P}$ и $(x_1, y_1), \dots, (x_n, y_n) \sim \mathcal{P}(du \times dv)$. Тогда можно записать две другие важные характеристики: ковариацию и коэффициент корреляции:

$$\begin{aligned} \text{cov}(\xi_1, \xi_2) &= \iint_{\mathbb{R}^2} (u - m_1(u))(v - m_1(v)) \mathcal{P}(du \times dv) & \widehat{\text{cov}}(\mathbf{x}, \mathbf{y}) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{\mathbf{x}})(y_i - \bar{\mathbf{y}}) \\ \text{cor}(\xi_1, \xi_2) &= \frac{\text{cov}(\xi_1, \xi_2)}{\sigma_{\xi_1} \sigma_{\xi_2}} & \widehat{\text{cor}}(\mathbf{x}, \mathbf{y}) &= \frac{\widehat{\text{cov}}(\mathbf{x}, \mathbf{y})}{s(\mathbf{x})s(\mathbf{y})}. \end{aligned}$$

Замечание (Важное). $\xi_1 \nparallel \xi_2 \implies \text{cov}(\xi_1, \xi_2) \neq 0$, но $\text{cov}(\xi_1, \xi_2) = 0 \not\implies \xi_1 \parallel \xi_2$. Необходимость и достаточность выполняется только в случае нормального распределения.

Замечание (Проблема моментов). Для заданной последовательности моментов m_1, m_2, \dots не обязательно существовать подходящее распределение. Помимо требований $m_{2k} \geq 0$ и взаимосвязи между соседними моментами по неравенству Гёльдера, существенно, что ряд Тейлора по m_ℓ , в который, как известно, раскладывается характеристическая функция, должен сходиться равномерно.

2.3. Свойства оценок

2.3.1. Несмещенность

Определение. Смещение³ есть

$$\text{bias } \hat{\theta}_n := E\hat{\theta}_n - \theta \quad \forall \theta \in \Theta.$$

Определение. Среднеквадратичная ошибка⁴ есть

$$\text{MSE } \hat{\theta}_n := E(\hat{\theta}_n - \theta)^2.$$

² «Skewness».

³ Bias.

⁴ Mean square error (MSE).

Замечание. Поскольку

$$D\hat{\theta}_n = D(\hat{\theta}_n - \theta) = E(\hat{\theta}_n - \theta)^2 - (E(\hat{\theta}_n - \theta))^2,$$

то

$$\underbrace{E(\hat{\theta}_n - \theta)^2}_{\text{MSE}} = D\hat{\theta}_n + \underbrace{(E(\hat{\theta}_n - \theta))^2}_{\text{bias}^2}.$$

Определение. Оценка называется *несмещенной*, если $\text{bias } \hat{\theta}_n = 0$, т.е.

$$E\hat{\theta}_n = \theta.$$

Пример. \bar{x} — несмещенная оценка $E\xi$.

Доказательство. Пусть $\theta = E\xi$, $\hat{\theta}_n = E\hat{\xi}_n = \bar{x}$. Тогда

$$E\bar{x} = E\frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{i=1}^n E x_i = \frac{1}{n} \sum_{i=1}^n E\xi = E\xi \implies E\hat{\theta}_n = E\theta, \text{ bias } \hat{\theta}_n = 0.$$

□

Пример. s^2 является только *асимптотически* несмещенной оценкой $D\xi$.

Доказательство. Поскольку дисперсия не зависит от сдвига, обозначим $\eta = \xi - E\xi$ и $y_i = x_i - E\xi$; тогда

$$\begin{aligned} E s^2 &= E\widehat{D\xi} = E\widehat{D\eta} = E\left(\widehat{E\eta^2} - \left(\widehat{E\eta}\right)^2\right) = E\left(\frac{1}{n} \sum_{i=1}^n y_i^2 - \left(\frac{1}{n} \sum_{i=1}^n y_i\right)^2\right) \\ &= E\left(\frac{1}{n} \sum_{i=1}^n y_i^2 - \frac{1}{n^2} \sum_{i,j=1}^n y_i y_j\right) = E\left(\frac{1}{n} \sum_{i=1}^n y_i^2 - \frac{1}{n^2} \sum_{i,j=1}^n y_i^2\right) = \frac{1}{n} \sum_{i=1}^n E y_i^2 - \frac{1}{n^2} \sum_{i=1}^n E y_i^2 \\ &= \frac{1}{n} \sum_{i=1}^n E (x_i - E\xi)^2 - \frac{1}{n^2} \sum_{i=1}^n E (x_i - E\xi)^2 = \frac{1}{n} \sum_{i=1}^n D x_i - \frac{1}{n^2} \sum_{i=1}^n D x_i = D\xi - \frac{1}{n} D\xi \\ &= \frac{n-1}{n} D\xi \xrightarrow{n \rightarrow \infty} D\xi. \end{aligned}$$

□

Определение. *Исправленная дисперсия:*

$$\tilde{s}^2 := \frac{n}{n-1} s^2.$$

2.3.2. Состоятельность

Определение. Оценка называется *состоятельной* в *среднеквадратичном смысле*, если

$$\text{MSE } \hat{\theta}_n \xrightarrow{n \rightarrow \infty} 0.$$

Определение. Оценка называется *состоятельной*, если

$$\hat{\theta}_n \xrightarrow{P} \theta.$$

Предложение. Если оценка несмещенная и состоятельная в среднеквадратичном смысле, то она состоятельная.

Доказательство. В самом деле, по неравенству Чебышева,

$$\mathbb{P}(|\hat{\theta}_n - \theta| > \epsilon) = \mathbb{P}(|\hat{\theta}_n - \mathbb{E}\hat{\theta}_n| > \epsilon) \leq \frac{D\hat{\theta}_n}{\epsilon^2} = \frac{\text{MSE } \hat{\theta}_n}{\epsilon^2} \xrightarrow{n \rightarrow \infty} 0.$$

□

Пример. $\hat{\mathbf{m}}_k$ является состоятельной оценкой \mathbf{m}_k .

Доказательство. Докажем для $\hat{\mathbf{m}}_1$. По определению выборки до эксперимента, $x_i \sim \mathcal{P}$. Тогда, по теореме Хинчина о ЗБЧ,

$$\bar{\mathbf{x}} = \frac{\sum_{i=1}^n x_i}{n} \xrightarrow{\mathbb{P}} \mathbf{m}_1(\mathcal{P}).$$

Для k -го момента доказываем аналогично заменой $y_i := x_i^k$.

□

Замечание. Для $\mathbf{m}_k^{(0)}$ доказательство не пройдет, потому что x_i и $\bar{\mathbf{x}}$ не будут независимыми.

Пример. $\hat{\mathbf{m}}_k^{(0)}$ является состоятельной оценкой $\mathbf{m}_k^{(0)}$.

Утверждение. Пусть $\xi_n \xrightarrow{\mathbb{P}} c$ и $f \in C(U_\epsilon(c))$. Тогда $f(\xi_n) \xrightarrow{\mathbb{P}} f(c)$.

Доказательство предложения. Докажем для s^2 . Пусть $f : (x, y) \mapsto x - y^2$. Устроим последовательность $(\hat{\mathbf{m}}_2, \hat{\mathbf{m}}_1) \xrightarrow{\mathbb{P}} (\mathbf{m}_2, \mathbf{m}_1)$. Тогда

$$f(\hat{\mathbf{m}}_2, \hat{\mathbf{m}}_1) = \hat{\mathbf{m}}_2 - \hat{\mathbf{m}}_1^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{\mathbf{x}}^2 = s^2 \xrightarrow{\mathbb{P}} f(\mathbf{m}_2, \mathbf{m}_1) = D\xi.$$

Для $\mathbf{m}_k^{(0)}$ доказываем аналогично.

□

Пример. $\bar{\mathbf{x}}$ — состоятельная оценка $\mathbb{E}\xi$.

Доказательство. Либо по (2.3.2) для $k = 1$, либо из того факта, что $\text{bias } \bar{\mathbf{x}} = 0$, значит

$$\text{MSE } \bar{\mathbf{x}} = D\bar{\mathbf{x}} = \frac{D\xi}{n} \xrightarrow{n \rightarrow \infty} 0,$$

и по (2.3.2) получаем утверждение.

□

Пример. s^2 — состоятельная оценка $D\xi$.

Доказательство. По (2.3.2) с $k = 2$.

□

Пример. $\widehat{\text{cdf}}_\xi$ — состоятельная оценка cdf в каждой точке.

Доказательство. Введем

$$y_i := \mathbf{1}_{\{x_i < x\}} = \begin{cases} 1 & x_i < x \\ 0 & x_i \geq x. \end{cases}$$

Тогда по теореме Хинчина о ЗБЧ,

$$\widehat{\text{cdf}}_\xi(x) = \frac{|\{x_i \in \bar{\mathbf{x}} : x_i < x\}|}{n} = \frac{\sum_{i=1}^n y_i}{n} \xrightarrow{\mathbb{P}} \mathbb{E}y_i = \mathbb{E}\mathbf{1}_{\{x_i < x\}} = \mathbb{P}(x_i < x) = \text{cdf}(x).$$

Независимость y_i очевидна, если расписать $\text{cdf}_{y_i}(x) \text{cdf}_{y_j}(y) = \text{cdf}_{y_i, y_j}(x, y)$.

□

Замечание. Помимо ранее упомянутых, больше нет состоятельных выборочных характеристик.

Утверждение. Пусть $\exists! p_0 : \text{cdf}(x) = p_0$ и $\text{cdf}(x)$ монотонно возрастает в окрестности p_0 . Тогда $\bar{z}_{p_0} \xrightarrow{\mathbb{P}} z_{p_0}$, т.е. является состоятельной оценкой.

2.3.3. Асимптотическая нормальность

Определение. Оценка называется *асимптотически нормальной*, если

$$\frac{\hat{\theta}_n - E\hat{\theta}_n}{\sqrt{D\hat{\theta}_n}} \xrightarrow{d} N(0, 1).$$

Замечание. Если $\xi \sim N(a, \sigma^2)$, то \bar{x} — просто нормальная оценка (как линейная комбинация нормальных случайных величин).

Все рассмотренные прежде оценки — асимптотически нормальные.

2.3.4. Эффективность

Определение. Соответствующая определению 2.5.2 оценка называется *эффективной*.

Определение. Оценка $\hat{\theta}^{(1)}$ называется *эффективной* в сравнении с $\hat{\theta}^{(2)}$, если

$$MSE \hat{\theta}^{(1)} < MSE \hat{\theta}^{(2)}.$$

Замечание. Для несмещенных оценок это эквивалентно, конечно,

$$D\hat{\theta}^{(1)} < D\hat{\theta}^{(2)}.$$

Пример (Сравнение оценок мат. ожидания симметричного распределения). Пусть \mathcal{P} *симметрично* — в этом случае $\widehat{\text{med}} \xi = \bar{x}$ и имеет смысл сравнить две этих характеристики.

$$\begin{aligned} D\bar{x} &= \frac{D\xi}{n} \\ \widehat{D\text{med}} \xi &\sim \frac{1}{4n \text{pdf}_{N(a, \sigma^2)}^2(\text{med } \xi)} \quad \text{при } n \rightarrow \infty. \end{aligned}$$

Так, если $\xi \sim N(a, \sigma^2)$, то

$$\text{pdf}_{N(a, \sigma^2)}^2(\text{med } \xi) = \frac{1}{2\pi\sigma^2} \exp \left\{ -\frac{(\text{med } \xi - a)^2}{\sigma^2} \right\} = \frac{1}{2\pi\sigma^2},$$

откуда

$$\widehat{D\text{med}} \xi = \frac{\pi \sigma^2}{2 n} > \frac{\sigma^2}{n} = D\bar{x},$$

значит $\widehat{\text{med}} \xi$ эффективнее \bar{x} .

Замечание. В то же время, $\widehat{\text{med}} \xi$ более устойчива к аутлаерам, чем \bar{x} , и этим лучше.

2.4. Метод моментов

Пусть $\mathcal{P}(\theta)$, $\theta = (\theta_1, \dots, \theta_r)^T$ — параметрическая модель. Найдём оценки для параметров $\hat{\theta}_i$, $i \in \overline{1:r}$, для чего составим и решим систему уравнений:

$$\begin{cases} m_1 = \phi_1(\theta_1, \dots, \theta_r) \\ \vdots \\ m_r = \phi_r(\theta_1, \dots, \theta_r) \end{cases} \implies \begin{cases} \theta_1 = f_1(m_1, \dots, m_r) \\ \vdots \\ \theta_r = f_r(m_1, \dots, m_r). \end{cases}$$

Примем

$$\hat{\theta}_i = f_i(\hat{m}_1, \dots, \hat{m}_r).$$

Замечание. Поскольку f_i — непрерывные функции и при непрерывных преобразованиях сходимость не портится, оценки $\hat{\theta}_i$ являются состоятельными.

Как правило, эти оценки смещенные.

Замечание. Случается, что решение находится вне пространства параметров. На практике, если пространство параметров компактное, можно взять точку, ближайшую к полученной оценке. Однако это свидетельствует о том, что модель плохо соответствует данным.

Пример 1 ($r = 1$). Пусть $\mathcal{P}_\xi(\lambda) = \text{Exp}(\lambda)$. Тогда $E\xi = 1/\lambda$ и $\bar{x} = 1/\lambda$.

Пример 2 ($r = 2$). Пусть $\mathcal{P}_\xi(\theta_1, \theta_2) = \text{Bin}(m, p)$. Тогда

$$\begin{cases} E\xi = mp \\ D\xi = mp(1-p) \end{cases} \quad \begin{cases} m = \frac{E\xi}{p} \\ D\xi = E\xi - E\xi p \end{cases} \quad \begin{cases} p = \frac{E\xi - D\xi}{E\xi} \\ m = \frac{(E\xi)^2}{E\xi - D\xi} \end{cases} \Rightarrow \begin{cases} \hat{p} = \frac{\bar{x} - s^2}{\bar{x}} \\ \hat{m} = \frac{\bar{x}^2}{\bar{x} - s^2}. \end{cases}$$

2.5. Метод оценки максимального правдоподобия

2.5.1. Функция правдоподобия и метод

Пусть $\mathcal{P}_\xi(\boldsymbol{\theta})$, $\boldsymbol{\theta} = (\theta_1, \dots, \theta_r)^T$ — параметрическая модель.

Определение. *Функция правдоподобия:*

$$\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{y}) = P(\mathbf{y} \mid \boldsymbol{\theta}) = \begin{cases} P_\boldsymbol{\theta}(x_1 = y_1, \dots, x_n = y_n) & \mathcal{P}_\xi(\boldsymbol{\theta}) \text{ дискретно;} \\ p_\boldsymbol{\theta}(\mathbf{y}) & \mathcal{P}_\xi(\boldsymbol{\theta}) \text{ абсолютно непрерывно.} \end{cases}$$

Пример 3. Пусть $\xi \sim N(a, \sigma^2)$. По независимости x_i , $p_\boldsymbol{\theta}(\mathbf{x})$ распадается в произведение:

$$\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}) = p_\boldsymbol{\theta}(\mathbf{x}) = \prod_{i=1}^n p_\boldsymbol{\theta}(x_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x_i - a)^2}{2\sigma^2}\right\} = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - a)^2\right\}.$$

Пример 4. $\xi \sim \text{Pois}(\lambda)$,

$$P(\xi = k) = \frac{\lambda^k}{k!} e^{-\lambda} \Rightarrow \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}) = \prod_{i=1}^n \frac{1}{x_i!} \lambda^{x_i} e^{-\lambda} = \frac{1}{\prod_{i=1}^n x_i!} \lambda^{n\bar{x}} e^{-n\lambda}.$$

Утверждение. Пусть \mathbf{x} — выборка. В качестве оценки максимального правдоподобия⁵ $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ следует взять

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \ln \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}).$$

Предложение. $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ является состоятельной оценкой.

Доказательство. Пусть $\boldsymbol{\theta}_0$ — истинный параметр $\mathcal{P}(\boldsymbol{\theta})$. По УЗБЧ,

$$\frac{1}{n} \ln \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \ln p_\boldsymbol{\theta}(x_i) \xrightarrow{P} E \ln p_\boldsymbol{\theta}(x_i) = \int_{\mathbb{R}} \ln(p_\boldsymbol{\theta}(x)) p_{\boldsymbol{\theta}_0}(x) dx.$$

Навесим на обе стороны argmax в условии, что это непрерывное преобразование:

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} \leftarrow \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \frac{1}{n} \ln \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}) \xrightarrow{P} \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \int_{\mathbb{R}} \ln(p_\boldsymbol{\theta}(x)) p_{\boldsymbol{\theta}_0}(x) dx \rightarrow \boldsymbol{\theta}^*.$$

Тогда в предположении непрерывности $p_\boldsymbol{\theta}$ по $\boldsymbol{\theta}$, $\hat{\boldsymbol{\theta}}_{\text{MLE}} \xrightarrow{P} \boldsymbol{\theta}^*$. Покажем, что $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$. Поделим на $p_{\boldsymbol{\theta}_0}$ — константу по $\boldsymbol{\theta}$:

$$\frac{1}{n} \sum_{i=1}^n \ln \frac{p_\boldsymbol{\theta}}{p_{\boldsymbol{\theta}_0}}(x_i) \xrightarrow{P} \int_{\mathbb{R}} \ln \left(\frac{p_\boldsymbol{\theta}(x)}{p_{\boldsymbol{\theta}_0}(x)} \right) p_{\boldsymbol{\theta}_0}(x) dx = E \ln \frac{p_\boldsymbol{\theta}}{p_{\boldsymbol{\theta}_0}} \leq \ln E \frac{p_\boldsymbol{\theta}}{p_{\boldsymbol{\theta}_0}} = \ln \int_{\mathbb{R}} \frac{p_\boldsymbol{\theta}}{p_{\boldsymbol{\theta}_0}}(x) p_{\boldsymbol{\theta}_0}(x) dx = \ln 1 = 0$$

⁵Maximum likelihood estimate (MLE).

по неравенству Ёнсена $\mathbb{E}g(\xi) \leq g(\mathbb{E}\xi)$, для выпуклой вверх $g(x) = \log(x)$. Таким образом,

$$\int_{\mathbb{R}} \ln \left(\frac{p_{\theta}(x)}{p_{\theta_0}(x)} \right) p_{\theta_0}(x) dx = 0 \iff \ln \left(\frac{p_{\theta}(x)}{p_{\theta_0}(x)} \right) = 0 \iff \frac{p_{\theta}(x)}{p_{\theta_0}(x)} = 1 \iff p_{\theta}(x) = p_{\theta_0}(x).$$

В предположении свойства *идентифицируемости* задачи $(\theta_1 \neq \theta_2 \implies \mathcal{P}_{\theta_1} \neq \mathcal{P}_{\theta_2})$, получаем $\theta = \theta_0$. \square

Пример. $\xi \sim \text{Pois}(\lambda)$.

$$\ln \mathcal{L}(\lambda \mid \mathbf{x}) = - \sum_{i=1}^n x_i! - n\lambda + \ln(\lambda^{n\bar{x}}) \implies \frac{d \ln \mathcal{L}(\lambda \mid \mathbf{x})}{d\lambda} = -n + \lambda^{-n\bar{x}} n\bar{x} \lambda^{n\bar{x}-1} = -n + \frac{n\bar{x}}{\lambda}$$

откуда

$$\frac{d \ln \mathcal{L}(\lambda \mid \mathbf{x})}{d\lambda} = 0 \iff -n + \frac{n\bar{x}}{\lambda} = 0, \quad n\bar{x} - n\lambda = 0, \quad \lambda = \bar{x}.$$

Утверждение. В условиях регулярности⁶:

1. Существует один глобальный максимум, так что

$$\left. \frac{d\mathcal{L}(\theta \mid \mathbf{x})}{d\theta} \right|_{\theta=\hat{\theta}_{\text{MLE}}} = 0.$$

2. $\hat{\theta}_{\text{MLE}}$ обладает всеми свойствами:

- a) Состоятельность;
- b) Асимптотическая несмещенность;
- c) Асимптотическая нормальность;
- d) Эффективность.

2.5.2. Информационное количество Фишера и эффективность оценки MLE

Пусть $r = 1$.

Определение. *Информанта n -го порядка:*

$$S_n(\mathbf{x}, \theta) = \frac{d^n \ln \mathcal{L}(\theta \mid \mathbf{x})}{d\theta^n}.$$

Определение. *Информационное количество Фишера:*

$$I_n(\theta) := -\mathbb{E} S_2(\mathbf{x}, \theta).$$

Утверждение.

$$I_n(\theta) = \mathbb{E} S_1^2(\mathbf{x}, \theta).$$

Пример. $\xi \sim \text{Pois}(\lambda)$.

$$S_1(\mathbf{x}, \theta) = -n + \frac{n\bar{x}}{\lambda}, \quad S_2(\mathbf{x}, \theta) = -\frac{n\bar{x}}{\lambda^2} \implies I_n(\lambda) = \mathbb{E} \frac{n\bar{x}}{\lambda^2} = \frac{n}{\lambda^2} \mathbb{E} \bar{x} = \frac{n}{\lambda}.$$

⁶Условия регулярности таковы:

- $\{x : p_{\theta}(x) > 0\}$ не зависит от θ ,
- $\int_{V^n} \hat{\theta} \mathcal{L}(\theta \mid \mathbf{x}) d\mathbf{x}$ и $\int_{V^n} \mathcal{L}(\theta \mid \mathbf{x}) d\mathbf{x}$ можно дифференцировать по θ под знаком интеграла,
- $I_n(\theta) > 0$.

Замечание.

$$\ln \mathcal{L}(\theta | \mathbf{x}) = \sum_{i=1}^n \ln p_{\theta}(x_i) \implies S_2 = \frac{d^2 \ln \mathcal{L}(\theta | \mathbf{x})}{d\theta^2} = \sum_{i=1}^n (\ln p_{\theta}(x_i))'',$$

откуда, для повторной независимой выборки,

$$I_n(\theta) = - \sum_{i=1}^n \mathbb{E}(\ln p_{\theta}(x_i))'' = n \cdot i(\theta), \quad \text{где } i(\theta) = -\mathbb{E}(\ln p_{\theta}(\xi))''.$$

Утверждение. Для несмещенных оценок в условиях регулярности справедливо неравенство Рао–Крамера:

$$D\hat{\theta}_n \geq \frac{1}{I_n(\theta)}.$$

Определение. Эффективная оценка:

$$D\hat{\theta}_n = \frac{1}{I_n(\theta)}.$$

Замечание. Для несмещенных оценок неравенство Рао–Крамера указывает точную нижнюю границу дисперсий оценок.

Пример. Пусть $\xi \sim \text{Pois}(\lambda)$. Поскольку

$$\begin{aligned} D\hat{\lambda}_n &= D\bar{x} = \mathbb{E}\xi/n = \lambda/n \\ I_n(\lambda) &= n/\lambda, \end{aligned}$$

то $\hat{\lambda}_n$ — эффективная оценка (как и ожидалась по свойствам $\hat{\theta}_{\text{MLE}}$).

Определение. Пусть $\hat{\theta}_n$ — асимптотически несмещенная оценка. Тогда $\hat{\theta}_n$ — асимптотически эффективная, если

$$D\hat{\theta}_n \cdot I_n \xrightarrow{n \rightarrow \infty} 1.$$

Пример. Пусть $\xi \sim N(a, \sigma^2)$. Можно посчитать, что s^2 является только асимптотически эффективной оценкой σ^2 ; \tilde{s}^2 — просто эффективной.

3. Некоторые распределения, связанные с нормальным

3.1. Распределение $\chi^2(m)$

Определение (Распределение $\chi^2(m)$). η имеет распределение χ^2 с m степенями свободы:

$$\eta \sim \chi^2(m) \iff \eta = \sum_{i=1}^m \zeta_i^2, \quad \zeta_i \sim N(0, 1), \quad \zeta_i \text{ независимы.}$$

Свойства⁷ $\chi^2(m)$

$$\begin{aligned} \mathbb{E}\eta &= \sum_{i=1}^m \mathbb{E}\beta_i^2 = m \\ D\eta &= 2m. \end{aligned}$$

Утверждение. Пусть $\eta_m \sim \chi^2(m)$. Тогда, по ЦПТ,

$$\frac{\eta_m - \mathbb{E}\eta_m}{\sqrt{D\eta_m}} = \frac{\eta_m - m}{\sqrt{2m}} \xrightarrow{d} N(0, 1).$$

Пример. $m = 50$, $\eta_m = 80$. Тогда

$$\frac{80 - 50}{10} = 3$$

и можно посчитать, к примеру, $\Phi(3) \approx 0.9986$.

⁷Вычисление $D\eta$: <https://www.statlect.com/probability-distributions/chi-square-distribution>

3.2. Распределение Стьюдента $t(m)$

Определение (Распределение $t(m)$). ξ имеет распределение Стьюдента с m степенями свободы, если

$$\xi \sim t(m) \iff \xi = \frac{\zeta}{\sqrt{\eta/m}}, \quad \zeta \sim N(0, 1), \quad \eta \sim \chi^2(m).$$

Свойства $t(m)$

- При $m = 1$ это распределение Коши.
- При $m > 1$, $E\zeta = 0$ по симметричности.
- При $m > 2$, $D\zeta = m/(m - 2)$.
- При $m > 3$, $A\zeta = 0$ по симметричности.
- При $m > 4$, $K\zeta = 6/(m - 4)$.

Предложение. Распределение Стьюдента сходится к стандартному нормальному:

$$t \Rightarrow N(0, 1).$$

Соображения по поводу. $D\zeta \rightarrow 1$, $K\zeta \rightarrow 0$. □

3.3. Распределение Фишера

Определение. Распределение Фишера имеет вид

$$F(m, k) = \frac{\chi^2(m)/m}{\chi^2(k)/k}.$$

Замечание. $F(1, k) \sim t^2(k)$; $m \cdot F(m, \infty) = \chi^2(m)$ потому что $\chi^2(k) \xrightarrow[k \rightarrow \infty]{} 1$.

3.4. Квадратичные формы от нормально распределенных случайных величин

Пусть $\xi = (\xi_1, \dots, \xi_p)^T \sim N(0, \sigma^2 I_p)$, B — симметричная, неотрицательно определенная матрица. Найдем распределение $\xi^T B \xi$.

Утверждение. Пусть $\xi \sim N(0, \sigma^2 I_p)$, B, C — симметричные матрицы размерности $p \times p$. Тогда $\xi^T B \xi \perp \xi^T C \xi \iff BC = 0$.

Пример (Независимость \bar{x}^2 и s^2). Пусть

$$B = \frac{1}{p} \begin{pmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{pmatrix};$$

тогда

$$\mathbf{x}^T B \mathbf{x} = \frac{1}{p} (x_1, \dots, x_p) \begin{pmatrix} \sum_{i=1}^p x_i \\ \vdots \\ \sum_{i=1}^p x_i \end{pmatrix} = \frac{1}{p} \sum_{j=1}^p \sum_{i=1}^p x_i x_j = \frac{1}{p} \left(\sum_{i=1}^p x_i \right)^2 = p \bar{x}^2.$$

В то же время, пусть

$$C = \begin{pmatrix} 1 - 1/p & \dots & -1/p \\ \vdots & \ddots & \vdots \\ -1/p & \dots & 1 - 1/p \end{pmatrix};$$

тогда

$$\begin{aligned} \mathbf{x}^T \mathbf{C} \mathbf{x} &= (x_1, \dots, x_p) \begin{pmatrix} (1 - 1/p) x_1 - 1/p \cdot \sum_{i=2}^p x_i \\ \vdots \\ -1/p \cdot \sum_{i=1}^{p-1} x_i + (1 - 1/p) x_p \end{pmatrix} = (x_1, \dots, x_p) \begin{pmatrix} x_1 - 1/p \cdot \sum_{i=1}^p x_i \\ \vdots \\ x_p - 1/p \cdot \sum_{i=1}^p x_i \end{pmatrix} \\ &= \sum_{j=1}^p \left(x_j^2 - \frac{1}{p} \sum_{i=1}^p x_i x_j \right) = \sum_{j=1}^p x_j^2 - p \bar{\mathbf{x}}^2 = p s^2. \end{aligned}$$

Но, учитывая $\mathbf{B}^2 = \mathbf{B}$ (легко проверяется),

$$\mathbf{C} = \mathbf{I}_p - \mathbf{B} \implies \mathbf{B} \mathbf{C} = \mathbf{B} - \mathbf{B}^2 = \mathbf{0},$$

откуда $\bar{\mathbf{x}}^2 \perp\!\!\!\perp s^2$.

Видно, что $\sigma^{-2} \boldsymbol{\xi}^T \mathbf{I}_p \boldsymbol{\xi} \sim \chi^2(p)$. На самом деле справедливо

Утверждение. Пусть $\boldsymbol{\xi} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_p)$, \mathbf{B} — симметричная, неотрицательно неопределенная матрица размерности $p \times p$ и $\text{rk } \mathbf{B} = r$. Тогда

$$\sigma^{-2} \boldsymbol{\xi}^T \mathbf{B} \boldsymbol{\xi} \sim \chi^2(r) \iff \mathbf{B}^2 = \mathbf{B}.$$

Пример. $n\sigma^{-2}s^2 \sim \chi^2(p-1)$. Воспользуемся представлением из предыдущего примера: $ps^2 = \mathbf{x}^T \mathbf{C} \mathbf{x}$. Но $\text{rk } \mathbf{C} = \text{rk}(\mathbf{I}_p - \mathbf{B}) = p-1$; $\mathbf{B}^2 = \mathbf{B}$, значит $p\sigma^{-2}s^2 \sim \chi^2(p-1)$.

Утверждение (Cochran). Пусть $\boldsymbol{\xi} \sim N(\mathbf{0}, \mathbf{I}_p)$, $\boldsymbol{\xi}^T \boldsymbol{\xi} = \sum_i Q_i$, где Q_i — квадратичная форма, заданная \mathbf{B}_i , $\text{rk } \mathbf{B}_i = r_i$. Тогда следующие утверждения эквивалентны:

1. $\sum r_i = p$
2. $Q_i \sim \chi^2(r_i)$
3. $Q_i \perp\!\!\!\perp Q_j, \quad \forall i \neq j$, т.е. $\mathbf{B}_i \mathbf{B}_j = \mathbf{0}$.

4. Проверка гипотез

Этот раздел иногда называется «Confirmatory Data Analysis» в противовес «Exploratory Data Analysis», не включающему в себя понятие *гипотезы*.

4.1. Построение критерия

4.1.1. Понятие гипотезы и критерия

Пусть $x_1, \dots, x_n \sim \mathcal{P}$, \mathcal{P} — множество всех распределений. О \mathcal{P} возможно делать утверждения вида $\mathcal{P} \in \mathcal{P}' \subset \mathcal{P}$. Стоит задача выбрать такое утверждение, что оно некоторым наилучшим образом соответствует выборке.

Определение. *Модель* — это предположение о выделенном классе $\mathcal{P}_M \subset \mathcal{P}$, которому принадлежит \mathcal{P} (допустим, $\mathcal{P}_M = \{N(\mu, \sigma_0)\}$, где σ_0 — фиксированное значение). Иными словами, это утверждение о \mathcal{P} , которое считается верным и не проверяется.

Определение. *Гипотеза* — утверждение о \mathcal{P} , требующее проверки. Гипотеза называется *простой*, если она соответствует только одному распределению в рамках рассматриваемой модели:

$$H : \mathcal{P} = \mathcal{P}_0 \in \mathcal{P}_M$$

(например, $\mathcal{P}_0 = N(\mu_0, \sigma_0)$) или *сложной*, если целому множеству:

$$H : \mathcal{P} \in \mathcal{P}' \subset \mathcal{P}_M$$

(например, $\mathcal{P}' = \{N(\mu, \sigma_0) : \mu > 0\}$).

Очень часто возникает (и далее рассматривается) случай выдвижения только лишь двух гипотез: $H_0 : \mathcal{P} \in \mathcal{P}_0 \subset \mathcal{P}$ — нулевой, основной и $H_1 : \mathcal{P} \in \mathcal{P}_1 \subset \mathcal{P}$ — альтернативной. H_1 учитывает отклонения от H_0 , обнаружение которых желательно. Возможны варианты:

- простая H_0 и простая H_1 ;
- простая H_0 и сложная H_1 ;
- сложная H_0 и сложная H_1 .

Определение. *Критерий* есть отображение

$$\varphi : \mathbf{x} \mapsto \{H_0, H_1\}.$$

Критерий «решает», противоречат или не противоречат выдвинутой гипотезе выборочные данные.

Определение. Говорят, что гипотеза *отвергается*, если $\varphi(\mathbf{x}) = H_1$ и не *отвергается* иначе.

Так как в заданной постановке любой критерий принимает не более двух значений, то $\text{dom } \varphi$ разбивается на два дизъюнктивных множества $\mathcal{A}_{\text{крит}}$ и $\mathcal{A}_{\text{дов}}$, называемых *критической* и *доверительной* областями, таких, что

$$\varphi(\mathbf{x}) = \begin{cases} H_0 & \mathbf{x} \in \mathcal{A}_{\text{дов}} \\ H_1 & \mathbf{x} \in \mathcal{A}_{\text{крит}}. \end{cases}$$

Поскольку выборка конечного объема позволяет делать только вероятностные заключения, со статистическим критерием ассоциированы ошибки i -ых родов.

Определение. Говорят, что произошла *ошибка i -го рода* критерия φ , если критерий отверг верную гипотезу H_{i-1} . Соответствующие вероятности обозначаются

$$\alpha_i(\varphi) = P_{H_{i-1}}(\varphi(\mathbf{x}) \neq H_{i-1}).$$

Поскольку в рассмотрение введены только H_0 и H_1 , возможны ошибки *I-го рода* — принятие случайного различия за систематическое и *II-го рода* — принятие наблюдаемого различия за случайный эффект с соответствующими вероятностями

$$\begin{aligned} \alpha_I &= P_{H_0}(\varphi(\mathbf{x}) \neq H_0) = P_{H_0}(\mathbf{x} \in \mathcal{A}_{\text{крит}}) \\ \alpha_{II} &= P_{H_1}(\varphi(\mathbf{x}) \neq H_1) = P_{H_1}(\mathbf{x} \in \mathcal{A}_{\text{дов}}). \end{aligned}$$

Замечание. Если H_i — сложная гипотеза, то $\alpha_{i+1}(\varphi)$ будет зависеть от того, на каком именно распределении \mathcal{P} , отвечающем H_i , вычисляется эта вероятность.

Определение. *Мощность* критерия против альтернативы это вероятность справедливо отвергнуть H_0 :

$$\beta = 1 - \alpha_{II} = 1 - P_{H_1}(\varphi(\mathbf{x}) = H_0) = P_{H_1}(\varphi(\mathbf{x}) = H_1).$$

Иными словами, это способность критерия отличать H_0 от H_1 .

Существует несколько подходов к сравнению критериев⁸. Выбор наиболее мощного критерия происходит так. До эксперимента фиксируют *уровень значимости*⁹ критерия $\alpha \in [0, 1]$ и рассматривают только критерии с $\alpha_I \leq \alpha$; среди них выбирают с наибольшей мощностью. Максимизировать β можно за счет правильного выбора $\mathcal{A}_{\text{крит}}$.

Замечание. Стандартные уровни значимости: $\alpha = 0.05$ или $\alpha = 0.01$.

⁸Минимаксный, байесовский, наиболее мощного критерия.

⁹Неформально, α обратно пропорциональна «строгости» критерия, выбираемой экспериментатором.

Определение. Критерий называется *состоятельным* против альтернативы H_1 , если $\forall P_1 \in \mathcal{P}_1$

$$\beta(\varphi, P_1) = 1 - P_{P_1}(\varphi(\mathbf{x}) = H_0) \xrightarrow{n \rightarrow \infty} 1.$$

Определение. Если

$\alpha = \alpha_I$ то критерий называется *точным*,

$\alpha \xrightarrow{n \rightarrow \infty} \alpha_I$ *асимптотическим*,

$\alpha < \alpha_I$ *радикальным* (т.е. отвергает гипотезу чаще, чем точный),

$\alpha > \alpha_I$ *консервативным* (если гипотеза отвергнута, то уж наверняка).

Замечание. Задача может допускать две постановки; в этом случае, поскольку α_I ($= \alpha$ для правильно построенного критерия) контролируется экспериментатором, проверяется отрицание эффекта, который хотят подтвердить: к примеру, что новое лекарство *не* лучше старого; если H_0 отвергнется, это будет означать, что новое лекарство-таки лучше старого с вероятностью не ниже, чем $1 - \alpha$.

Пример (С гранатами). FIXME

Замечание. Утверждать об *отвержении* гипотезы можно с вероятностью ошибки α (достаточно малой и произвольно задаваемой экспериментатором); утверждать о *принятии* гипотезы можно с вероятностью ошибки α_{II} — не контролируемой и потенциально довольно большой. Иными словами, попадание в доверительную область может означать как то, что H_0 верна, так и то, что верна H_1 , но для распознавания этого не хватило мощности. Поэтому безопасно гипотезу можно только отвергать или не отвергать. Можно и принять, если известна мощность критерия против всех возможных альтернатив, экспериментатора устраивающая.

Замечание. При высокой вероятности ошибки II-го рода возможна ситуация не отвержении заведомо ложной гипотезы. Это, в свою очередь, может произойти из-за маленького объема выборки (критерий не находит разницу, см. 4.1.2). Чем больше объем выборки, тем мощность больше, но возможна ситуация, когда критерий чувствителен настолько, что находит разницу там, где не должен — например, при генерации «идеальным» датчиком случайных чисел, начиная с какого-то объема заведомо истинная гипотеза может быть отвергнута из-за ошибок в точности представления чисел с плавающей точкой.

Определение. Критерий называется *одно- (двух-) сторонним* по тому, где находится альтернатива.

Определение. Критическая область называется *одно- (дву-) сторонней* по тому, где формально располагается $\mathcal{A}_{\text{крит}}$.

4.1.2. Построение критерия при помощи статистики критерия

Определение. Статистика критерия есть отображение

$$T : \mathbf{x} \mapsto y \in \mathbb{R}$$

такое, что при верной H_0 , $T \stackrel{d}{\rightarrow} Q$, где Q — полностью известное непрерывное распределение, а при верной H_1 известно поведение T .

Поскольку распределение T при верной H_0 известно, она должна вести себя как любая другая случайная величина из Q — попадание в некоторые области менее вероятно, чем в другие. Поэтому разумно разбить $\text{im } T$ по уровню значимости α на $\mathcal{A}_{\text{крит}} \sqcup \mathcal{A}_{\text{дов}}$ так, что попадание в $\mathcal{A}_{\text{крит}}$ происходит с заранее зафиксированной (малой) вероятностью α . Значит, если $T(\mathbf{x}) \in \mathcal{A}_{\text{крит}}$, то с

некоторой же вероятностью можно заявлять об отвержении H_0 . Таким образом, T измеряет то, насколько выборка соответствует гипотезе.

Разберем на примере построение разбиения. Пусть $\xi \sim N(\mu, \sigma^2)$, $H_0 : \mu = \mu_0$ и фиксирован α . По ??, используется статистика

$$T = z = \sqrt{n} \frac{\bar{x} - a_0}{\sigma} \sim N(0, 1).$$

В зависимости от H_1 , возможны варианты.

Простая альтернатива Пусть $H_1 : \mu = \mu_1$, причем $\mu_1 > \mu_0$. Тогда, поскольку при верной H_1 , $E\bar{x} = 1/n \cdot \sum_{i=1}^n \xi_i = n/n \cdot \mu_1$, то

$$ET = \frac{\sqrt{n}(\mu_1 - \mu_0)}{\sigma} \implies T \sim N\left(\frac{\sqrt{n}(\mu_1 - \mu_0)}{\sigma}, 1\right) \text{ при верной } H_1.$$

(дисперсия, конечно, не меняется при сдвиге).

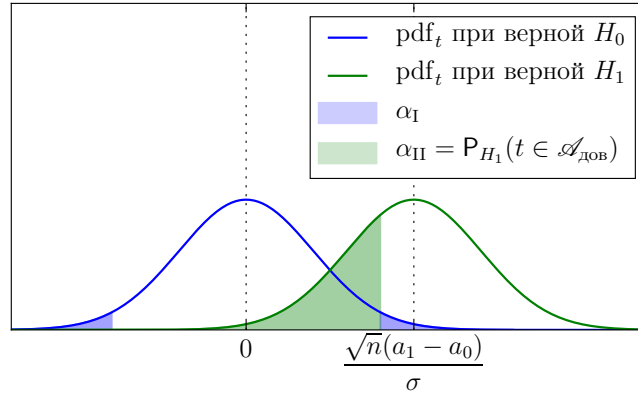


Рис. 1: Плотности распределения z (неоптимальное разбиение)

Чтобы минимизировать α_{II} , логично определить $\mathcal{A}_{\text{крит}}$ только на одном хвосте — с той стороны, где находится альтернатива. Помимо этого, по рисунку видно, что минимизировать α_{II} (согласившись на бóльшую ошибку первого рода) можно сдвинув вправо центр второй плотности, увеличив n . Аналогично, чем μ_1 дальше от μ_0 , тем α_{II} меньше.

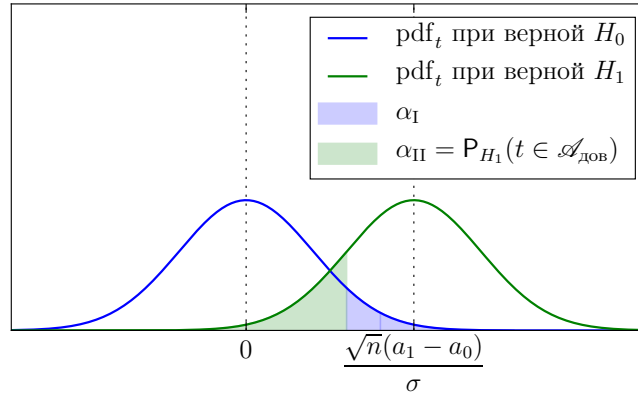


Рис. 2: Плотности распределения z (оптимальное разбиение)

Таким образом, $\mathcal{A}_{\text{крит}} = (C, +\infty)$, $C = z_{1-\alpha}$.

Разумеется, если $\mu_1 < \mu_0$, то $\mathcal{A}_{\text{крит}} = (-\infty, C)$, $C = z_\alpha$.

Односторонний критерий (сложная альтернатива) В общем случае, пусть $H_1 : \mu = \mu_1 \forall \mu_1 > \mu_0$; тогда по ЗБЧ $\bar{x} - \mu_0 \rightarrow \mu_1 - \mu_0 > 0$ и $T \rightarrow +\infty$. Следовательно, чтобы максимизировать величину $\beta = P_{H_1}(\mathbf{x} \in \mathcal{A}_{\text{крит}}^{(\alpha)})$, следует разместить \mathcal{A} на правом хвосте плотности: $\mathcal{A}_{\text{крит}} = (z_{1-\alpha}, \infty)$. Для $H_1 : \mu = \mu_1 \forall \mu_1 < \mu_0$ аналогично $\mathcal{A}_{\text{крит}} = (-\infty, z_\alpha)$.

Двусторонний критерий Пусть $H_1 : E\xi = \mu_1 \neq \mu_0$; тогда по ЗБЧ $\bar{x} - \mu_0 \rightarrow \mu_1 - \mu_0$ и $|T| \xrightarrow{n \rightarrow \infty} \infty$, откуда $\mathcal{A}_{\text{крит}} = \mathbb{R} \setminus (z_{\alpha/2}, z_{1-\alpha/2}) = \mathbb{R} \setminus (-C, C)$, где $C = -z_{\alpha/2}$.

В общем виде, с использованием статистики, критерий может быть определен как

$$\varphi(\mathbf{x}) = \begin{cases} H_0 & |T(\mathbf{x})| < C \\ H_1 & |T(\mathbf{x})| \geq C, \end{cases}$$

где критическое значение C определяется из уравнения $\alpha = P(|T| \geq C)$.

Иногда вместо сравнения значения T с критическим вычисляют *реально достигнутый уровень значимости критерия* («*p-value*»).

Определение. *p-value* значения статистики T на выборке \mathbf{x} есть вероятность, взяв выборку из распределения H_0 , получить по ней большее отклонение $|T(\mathbf{x})|$ эмпирического от истинного распределения, чем получено по проверяемой выборке:

$$p\text{-value} = \alpha^* = P_{H_0}(|T| \geq |T(\mathbf{x})|).$$

Значит, критерий может быть задан и как

$$\varphi(\mathbf{x}) = \begin{cases} H_0 & \alpha^* > \alpha \\ H_1 & \alpha^* \leq \alpha. \end{cases}$$

Замечание. *p-value* обратно пропорционален «существенности» результата.

Замечание. Пусть $\alpha^* = 0.05$. Это значит, что в среднем всего лишь 5% «контрольных» выборок, удовлетворяющих основной гипотезе, будут обладать большим отклонением $|T(\mathbf{x})|$ по сравнению с тестируемой выборкой — последняя ведет себя не хуже, чем 5% «правильных» выборок.

4.1.3. Схема построения критерия с помощью статистики

1. Фиксируют предположение относительно данных.
2. Выдвигают H_0 и H_1 .
 - H_0 формулируется согласно замечанию 4.1.1.
 - H_1 ставится по смыслу задачи (см. далее).
3. Выбирают подходящий критерий и статистику T .
4. Фиксируют уровень значимости α .
5. Строят разбиение $\text{im } T$ с помощью квантилей распределения T (при верной H_0) так, чтобы $\alpha_1 = \alpha$; положение квантилей выбирают из известного поведения статистики при верной H_1 .
6. Считают значение статистики и принимают решение об отвержении H_0 одним из способов:

$$\varphi(\mathbf{x}) = \begin{cases} H_0 & |T(\mathbf{x})| < C \\ H_1 & |T(\mathbf{x})| \geq C, \end{cases} \quad \varphi(\mathbf{x}) = \begin{cases} H_0 & \alpha^* > \alpha \\ H_1 & \alpha^* \leq \alpha. \end{cases}$$

Пример (Средняя температура в холодильнике). Хотят купить холодильник, такой, чтобы температура держалась в окрестности 0. Известно количество измерений $n = 25$ и $\bar{x} = 0.7$.

1. Пусть $\xi \sim N(a, 4)$.
2. Выдвинута $H_0 : E\xi = a_0 = 0$ — если гипотеза опровергнется, то холодильник не купят;
 $H_1 : E\xi = a_1 \neq a_0$.
3. Поскольку модель нормальная и известная σ^2 , выберем статистику ?? («z-test»):

$$z = \frac{\sqrt{n}(\bar{x} - a_0)}{\sigma} \sim N(0, 1) \text{ при верной } H_0.$$

Идеальное значение статистики — 0.

4. Зафиксируем два уровня значимости: $\alpha^{(1)} = 0.2$ (храним петрушку) и $\alpha^{(2)} = 0.01$ (храним дорогую красную икру).
5. Построим разбиение. Поскольку $a_1 \neq a_0$, то $\mathcal{A}_{\text{крит}} = \mathbb{R} \setminus (z_{\alpha/2}, z_{1-\alpha/2})$. Для введенных уровней значимости это означает

- a) $\mathcal{A}_{\text{крит}}^{(\alpha^{(1)})} \approx \mathbb{R} \setminus (-1.28, 1.28)$.
- b) $\mathcal{A}_{\text{крит}}^{(\alpha^{(2)})} \approx \mathbb{R} \setminus (-2.576, 2.576)$.

6. Посчитаем

$$z(\mathbf{x}) = \frac{\sqrt{n}(\bar{x} - a_0)}{\sigma} = \frac{5(0.7 - 0)}{2} = 1.75.$$

Дальнейшее принятие решения возможно на основании критического значения или p -value.

- По вычислению критического значения:

◇ $z \in \mathcal{A}_{\text{крит}}^{(\alpha^{(1)})}$, H_0 отвергается, холодильник не покупают.

◇ $z \in \mathcal{A}_{\text{дов}}^{(\alpha^{(2)})}$, H_0 не отвергается, холодильник, быть может, покупают.

- Можно посчитать p -value:

$$2 \cdot (1 - \text{cdf}_{N(0,1)}(1.75)) \approx 0.08.$$

Поэтому при уровне значимости $\alpha^{(1)} = 0.2 > 0.08$ H_0 отвергается, а при $\alpha^{(2)} = 0.01 < 0.08$ не отвергается.

Пример (С мышой). В одном из рукавов Т-образного лабиринта лежит морковка. К развилке по лабиринту бежит мышь и 7 раз из 10 поворачивает в направлении морковки. На основании этих данных хотим сделать вывод, что мышь чует морковь на расстоянии, после чего написать научную статью.

- $\xi \sim \text{Ber}(p)$. Выдвинем гипотезу, что мышь *не* чует морковку, $H_0 : p = p_0 = 0.5$. Поскольку $E\xi = p$, воспользуемся критерием для проверки гипотезы о значении среднего с идеальным значением 0; учитывая $D\xi = p(1 - p)$,

$$\begin{aligned} T &= \frac{\sqrt{n}(\bar{x} - p_0)}{\sqrt{p_0(1 - p_0)}} \xrightarrow{d} N(0, 1). \\ &= \frac{\sqrt{10} \cdot 0.2}{0.5} \approx 1.2649 \implies p\text{-value} = 2 \cdot (1 - \text{cdf}_{N(0,1)}(1.2649)) \approx 0.2. \end{aligned}$$

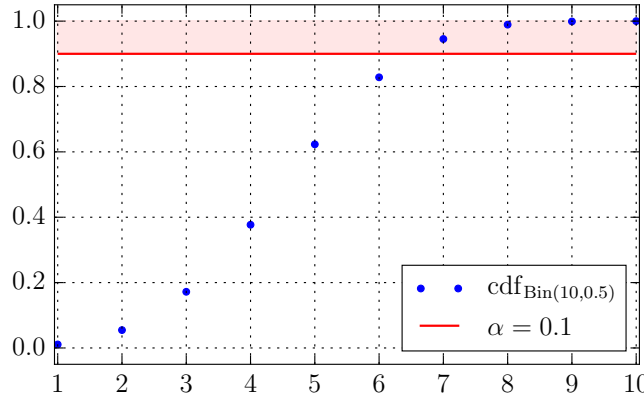
Значит, с уровнем значимости 0.2 гипотеза не отвергается. Хочется иметь, конечно, один из стандартных уровней значимости, например 0.1.

- Увеличим мощность критерия, введя альтернативную гипотезу, что мышь чует морковку (в предположении, что все мыши любят морковь и к ней бегут), $H_1 : p_1 > p_0$. По 4.1.2, можем устроить односторонний критерий, так что p -value теперь 0.1. Однако пользуемся асимптотическим критерием при $n = 10$.

- Воспользуемся точным односторонним критерием со статистикой

$$T := n\bar{x} = \sum_{i=1}^n x_i \sim \text{Bin}(n, p_0)$$

и идеальным значением np_0 . Тогда $T = 10 \cdot 0.7 = 7$. При уровне значимости $\alpha = 0.1$ успешно попадаем в критическую область, вследствие чего H_0 отвергается, и можем публиковаться.



Замечание. Исторически существовало два подхода к проверке гипотез: Фишера («significance test») и Неймана-Пирсона («hypothesis testing»).

Фишер Выдвигается H_0 . Подсчитывается и сообщается точное p -value. Если результат «незначительный», не делается никаких выводов об отвержении H_0 , но делается возможным дополнительный сбор данных.

Нейман-Пирсон Выдвигаются H_1, H_2 , фиксируются α_I, α_{II} и n . На этом основании определяются $\mathcal{A}_{\text{крит}}$ для каждой гипотезы. Если данные попали в $\mathcal{A}_{\text{крит}}$ H_1 — предпочитается H_2 , иначе H_2 .

Современная теория проверки гипотез есть смесь двух этих подходов, не всегда консистентная. Вводные курсы по статистике формулируют теорию, похожую на significance testing Фишера; при повышенных требованиях к математической строгости, пользуются теорией Неймана-Пирсона.

Замечание (О графике p -values). Поскольку

$$\alpha_I \leftarrow P_{H_0}(T \in \mathcal{A}_{\text{крит}}) = P_{H_0}(p\text{-value} < \alpha),$$

то p -value по распределению стремятся к $U(0,1)$ при верной H_0 . Это соображение позволяет визуально проверить истинность гипотезы: достаточно несколько (много) раз произвести эксперимент, для каждой выборки $\bar{x}^{(i)}$ посчитать свой p -value, построить график и убедиться, что получилась прямая. Для подсчета мощности $\beta = P_{H_1}(\mathbf{x} \in \mathcal{A}_{\text{крит}}^{(\alpha)}) = P_{H_1}(p\text{-value} < \alpha)$ считать выборку с параметрами H_1 , а T относительно H_0 .

4.2. Проверка гипотезы о значении мат. ожидания (t -критерий)

$H_0 : E\xi = a = a_0$. Соответствие оценки математического ожидания гипотезе удобно выражать разницей $\bar{x} - a_0$ с «идеальным» значением 0. Отнормировав эту разницу, получим статистику, распределение которой известно.

4.2.1. $D\xi = \sigma^2 < \infty$

Предложение. Пусть $D\xi = \sigma^2 < \infty$; тогда используется следующая статистика

$$t = \sqrt{n} \frac{(\bar{x} - a_0)}{\sigma} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

Доказательство. По ЦПТ. □

Предложение. При условии нормальности данных, t -критерий называется « z -критерием», причем

$$t = z \sim N(0, 1).$$

Доказательство.

$$z = \frac{\bar{x} - a_0}{\sqrt{D\bar{x}}} = \sqrt{n} \frac{\bar{x} - a_0}{\sigma} \sim N(0, 1).$$

□

Разбиение

$$H_1 : E\xi \neq a_0 \quad \mathcal{A}_{\text{крит}} = \mathbb{R} \setminus (z_{\alpha/2}, z_{1-\alpha/2})$$

$$H_1 : E\xi > a_0 \quad \mathcal{A}_{\text{крит}} = (z_{1-\alpha}, \infty)$$

$$H_1 : E\xi < a_0 \quad \mathcal{A}_{\text{крит}} = (-\infty, z_{\alpha})$$

4.2.2. $D\xi$ неизвестна

Предложение. Пусть $D\xi$ неизвестна; тогда используется следующая статистика

$$t = \sqrt{n-1} \frac{\bar{x} - a_0}{s} = \frac{\sqrt{n-1}(\bar{x} - a_0)}{\sqrt{n-1}/\sqrt{n} \cdot \tilde{s}} = \sqrt{n} \frac{\bar{x} - a_0}{\tilde{s}} \xrightarrow{n \rightarrow \infty} N(0, 1).$$

Предложение. При условии нормальности данных,

$$t \sim t(n-1).$$

Доказательство.

$$t = \frac{\sqrt{n-1}(\bar{x} - a_0)}{s} = \frac{\sqrt{n-1} \left(\frac{\bar{x} - a_0}{\sigma} \right)}{s/\sigma} = \frac{\left(\frac{\bar{x} - a_0}{\sigma} \right)}{\sqrt{\frac{s^2/\sigma^2}{n-1}}} = \frac{\frac{\sqrt{n}(\bar{x} - a_0)}{\sigma}}{\sqrt{\frac{ns^2/\sigma^2}{n-1}}} = \frac{\beta}{\sqrt{\eta/(n-1)}} \sim t(n-1),$$

поскольку

$$\beta = \frac{\sqrt{n}(\bar{x} - a_0)}{\sigma} \sim N(0, 1), \quad \eta = \frac{ns^2}{\sigma^2} \sim \chi^2(n-1).$$

□

Разбиение

$$H_1 : E\xi \neq a_0 \quad \mathcal{A}_{\text{крит}} = \mathbb{R} \setminus (\text{qnt}_{t(n-1)}(\alpha/2), \text{qnt}_{t(n-1)}(1 - \alpha/2))$$

$$H_1 : E\xi > a_0 \quad \mathcal{A}_{\text{крит}} = (\text{qnt}_{t(n-1)}(1 - \alpha), \infty)$$

$$H_1 : E\xi < a_0 \quad \mathcal{A}_{\text{крит}} = (-\infty, \text{qnt}_{t(n-1)} \alpha)$$

Замечание. При нормальной аппроксимации $\text{qnt}_{t(n-1)}$ заменить на $N(0, 1)$.

z -критерий для пропорции в модели Бернулли Пусть $\xi \sim \text{Ber}(p)$. Поскольку $E\xi = p$, можно воспользоваться только что введенной статистикой; учитывая $D\xi = p(1-p)$,

$$T = \sqrt{n} \frac{\bar{x} - p_0}{\sqrt{p_0(1-p_0)}} \xrightarrow{d} N(0, 1).$$

Разбиение будет таким же, как и в случае известной дисперсии.

4.3. Проверка гипотезы о значении дисперсии в нормальной модели (критерий χ^2)

Пусть $\xi \sim N(a, \sigma^2)$. $H_0 : D\xi = \sigma^2 = \sigma_0^2$. Соответствие оценки дисперсии гипотезе удобно выражать отношением s^2/σ_0^2 (или s_a/σ_0^2 если a известно) с «идеальным» значением 1. Домножив на n , получим статистику, распределение которой известно.

Замечание (Важное). Критерий работает только в нормальной модели и не становится асимптотически нормальным в ином случае!

4.3.1. $E\xi = a < \infty$

Предложение. Пусть $E\xi = a < \infty$; При условии нормальности данных используется следующая статистика:

$$\chi^2 = n \frac{s_a^2}{\sigma_0^2} \sim \chi^2(n).$$

Доказательство.

$$\chi^2 = \frac{ns_a^2}{\sigma_0^2} = \frac{n \cdot 1/n \cdot \sum_{i=1}^n (x_i - a)^2}{\sigma_0^2} = \sum_{i=1}^n \left(\frac{x_i - a}{\sigma_0} \right)^2 \sim \chi^2(n).$$

□

Разбиение

$$H_1 : D\xi \neq \sigma_0^2 \quad \mathcal{A}_{\text{крит}} = \mathbb{R}_+ \setminus \left(\text{qnt}_{\chi^2(n)}(\alpha/2), \text{qnt}_{\chi^2(n)}(1 - \alpha/2) \right)$$

$$H_1 : E\xi > a \quad \mathcal{A}_{\text{крит}} = (\text{qnt}_{\chi^2(n)}(1 - \alpha), \infty)$$

$$H_1 : E\xi < a \quad \mathcal{A}_{\text{крит}} = (0, \text{qnt}_{\chi^2(n)} \alpha)$$

4.3.2. $E\xi$ неизвестно

Предложение. Пусть $E\xi$ неизвестно. При условии нормальности данных используется следующая статистика:

$$\chi^2 = n \frac{s^2}{\sigma_0^2} = (n-1) \frac{\tilde{s}^2}{\sigma_0^2} \sim \chi^2(n-1).$$

Доказательство. См. (3.4).

□

Альтернативное доказательство. По определению запишем

$$D\hat{\xi}_n = D(\hat{\xi}_n - a) = E \left(\hat{\xi}_n - a \right)^2 - (E \left(\hat{\xi}_n - a \right))^2.$$

Но

$$\begin{aligned} D\hat{\xi}_n &= E(\hat{\xi}_n - E\hat{\xi}_n)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = s^2 \\ E(\hat{\xi}_n - a)^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - a)^2 = s_a^2 \\ (E(\hat{\xi}_n - a))^2 &= (\bar{x} - a)^2, \end{aligned}$$

откуда

$$s^2 = s_a^2 - (\bar{x} - a)^2.$$

Домножив обе части на n/σ_0^2 , получим

$$\frac{ns^2}{\sigma_0^2} = \frac{ns_a^2}{\sigma_0^2} - \frac{n(\bar{\mathbf{x}} - a)^2}{\sigma_0^2} = \underbrace{\frac{ns_a^2}{\sigma_0^2}}_{\sim \chi^2(n)} - \underbrace{\left(\frac{\sqrt{n}(\bar{\mathbf{x}} - a)}{\sigma_0} \right)^2}_{\sim \chi^2(1)} \Rightarrow \frac{ns^2}{\sigma_0^2} \sim \chi^2(n-1).$$

□

Замечание. Для строгого доказательства, нужно использовать независимость $\bar{\mathbf{x}}^2$ и s^2 (см. 3.4).

Разбиение

$$H_1 : D\xi \neq \sigma_0^2 \quad \mathcal{A}_{\text{крит}} = \mathbb{R}_+ \setminus \left(\text{qnt}_{\chi^2(n-1)}(\alpha/2), \text{qnt}_{\chi^2(n-1)}(1 - \alpha/2) \right)$$

$$H_1 : E\xi > a \quad \mathcal{A}_{\text{крит}} = (\text{qnt}_{\chi^2(n-1)}(1 - \alpha), \infty)$$

$$H_1 : E\xi < a \quad \mathcal{A}_{\text{крит}} = (0, \text{qnt}_{\chi^2(n-1)} \alpha)$$

Упражнение. $s^2 = 1.44$, $\bar{\mathbf{x}} = 55$, $n = 101$. Проверить гипотезу $\sigma_0^2 = 1.5$ в нормальной модели.

Решение. Воспользуемся статистикой

$$\chi^2 = \frac{ns^2}{\sigma_0^2} = 101 \cdot 0.96 = 96.96.$$

«Идеальные» значения близки к $E\xi_{\chi^2(100)} = 100$, так что определим критическую область на концах плотности:

$$p\text{-value}/2 = \text{cdf}_{\chi^2(100)}(96.96) = \text{pchisq}(96.96, 100) \approx 0.43 \Rightarrow p\text{-value} \approx 0.86.$$

Замечание. Можно посчитать и по таблицам для нормального распределения. Раз

$$\frac{\eta_m - E\eta_m}{\sqrt{D\eta_m}} \xrightarrow[m \rightarrow \infty]{d} N(0, 1),$$

то

$$\frac{96.96 - 100}{\sqrt{200}} \approx -0.215 \Rightarrow p\text{-value}/2 = \Phi(-0.215) \approx 0.415.$$

┘

4.4. Критерий χ^2 согласия с видом распределения

По выборке возможно проверить гипотезу о виде распределения случайной величины, реализацией которой является выборка.

Утверждение. Для проверки гипотезы согласия с видом произвольного *дискретного* распределения используется асимптотический критерий χ^2 («chi-squared test for goodness of fit»).

4.4.1. Распределение с известными параметрами

Пусть

$$H_0 : \mathcal{P} = \mathcal{P}_0, \text{ где } \mathcal{P}_0 : \begin{pmatrix} x_1^* & \dots & x_k^* \\ p_1 & \dots & p_k \end{pmatrix}.$$

Сгруппируем \mathbf{x} ; каждому x_i^* сопоставим *эмпирическую* абсолютную частоту ν_i ; тогда np_i — *ожидаемая* абсолютная частота.

В качестве меры расхождения между эмпирическим и генеральным распределением рассматривается величина

$$\sum_{i=1}^k c_i \left(\frac{\nu_i}{n} - p_i \right)^2, \quad c_i = \frac{n}{p_i},$$

откуда записывается статистика критерия

$$T = \sum_{i=1}^k \frac{(\nu_i - np_i)^2}{np_i}$$

с идеальным значением 0.

Утверждение. $T \xrightarrow{d} \chi^2(k-1)$.

Разбиение $\mathcal{A}_{\text{крит}} = (\text{qnt}_{\chi^2(k-1)}(1-\alpha), \infty)$ — гипотеза отвергается, если расстояние между предполагаемым и наблюдаемым распределениями большое.

Упражнение. $n = 100$,

$$\begin{pmatrix} \diamond & \heartsuit & \clubsuit & \spadesuit \\ 20 & 30 & 10 & 40 \end{pmatrix}.$$

Проверить гипотезу, что колода полная.

Решение. $H_0 : \mathcal{P}_\xi = \text{U}(1/4)$. Поскольку речь идет о согласии с дискретным не параметризованным распределением, напрямую воспользуемся критерием χ^2 . Раз все $np_i = 100 \cdot 1/4 = 25$,

$$\chi^2 = \sum_{i=1}^k \frac{(\nu_i - np_i)^2}{np_i} = 1 + 1 + \frac{15^2}{25} + \frac{15^2}{25} = 2 + 2 \cdot 9 = 20.$$

Так как $\chi^2 \sim \chi^2(k-1) = \chi^2(3)$ со средним 3, и «идеальное» значение 0, определим критическую область в правом конце плотности. Из этих соображений

$$p\text{-value} = 1 - \text{cdf}_{\chi^2(3)}(20) = 1 - \text{pchisq}(20, 3) \approx 0.00017.$$

┘

4.4.2. Распределение с неизвестными параметрами

В случае сложной гипотезы $\mathcal{P} \in \{\mathcal{P}(\boldsymbol{\theta})\}_{\boldsymbol{\theta} \in \Theta}$, $\boldsymbol{\theta} = (\theta_1, \dots, \theta_r)^T$, следует найти оценку $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ (или $\hat{\boldsymbol{\theta}} : \hat{\boldsymbol{\theta}} \rightarrow \hat{\boldsymbol{\theta}}_{\text{MLE}}$) по методу максимального правдоподобия. При подстановке оценок вместо истинных параметров критерий становится консервативным. Чтобы этого избежать, необходимо сделать поправку на количество параметров — отнять r . Что приятно, одна и та же поправка работает для всех распределений; в этом случае,

$$T \xrightarrow{d} \chi^2(k - r - 1).$$

Упражнение 1. 60 человек купило подарок сразу, 10 со второго раза, 20 с третьего, 10 с четвертого:

$$\begin{pmatrix} 0 & 1 & 2 & 3 \\ 60 & 10 & 20 & 10 \end{pmatrix}.$$

Проверить гипотезу о том, что это выборка из геометрического распределения.

Решение. $H_0 : \mathcal{P}_\xi = \text{Geom}(p)$. Воспользуемся критерием χ^2 для параметризованного распределения $\text{Geom}(\hat{p}_{\text{MLE}})$.

Найдем

$$\hat{p}_{\text{MLE}} = \underset{p}{\operatorname{argmax}} \log \mathcal{L}(\mathbf{x}; p) \iff \frac{d}{dp} \log \mathcal{L}(\mathbf{x}; \hat{p}_{\text{MLE}}) = 0.$$

Так как $\text{pdf}_{\text{Geom}(p)}(k) = (1-p)^{k-1}p$,

$$\begin{aligned} \log \mathcal{L}(\mathbf{x}; p) &= \log \prod_{k=1}^n (1-p)^{k-1} p = \log(1-p)^{n\bar{x}} p^n = n\bar{x} \log(1-p) + n \log p \\ &= n(\bar{x} \log(1-p) + \log p) \end{aligned}$$

откуда

$$\frac{d}{dp} \log \mathcal{L}(\mathbf{x}; p) = n \left(-\frac{\bar{\mathbf{x}}}{1-p} + \frac{1}{p} \right) = 0 \iff 1 - p - p\bar{\mathbf{x}} = 0 \iff p = \frac{1}{1 + \bar{\mathbf{x}}}.$$

Учитывая

$$\bar{\mathbf{x}} = 0.1 + 2 \cdot 0.2 + 3 \cdot 0.1 = 0.8,$$

найдем

$$\hat{p}_{\text{MLE}} = \frac{1}{1 + 0.8} \approx 0.55.$$

Посчитаем статистику χ^2 , найдя соответствующие p_i :

$$p_0 = P_{\text{Geom}(0.55)}(0) = 0.55, \quad p_1 \approx 0.26, \quad p_2 \approx 0.11, \quad p_3 \approx 0.09.$$

Тогда

$$\chi^2 = \sum_{i=1}^k \frac{(\nu_i - np_i)^2}{np_i} = \frac{25}{55} + \frac{16^2}{26} + \frac{81}{11} + \frac{1}{9} \approx 17.77.$$

Наконец, поскольку $\chi^2 \xrightarrow{n \rightarrow \infty} \chi^2(k - r - 1)$,

$$p\text{-value} = 1 - \text{cdf}_{\chi^2(2)}(17.77) \approx 0.00014.$$

┘

Определение. Критерий *применим*, если $\alpha \rightarrow \alpha_I$.

Замечание. Поскольку критерий асимптотический, с достаточной степенью точностью он применим в случае, если

1. $n \geq 50$;
2. $np_i \geq 5$.

Замечание. Если условие $np_i \geq 5$ не выполняется, следует объединить состояния, например, с краев или слева направо; если в хвосте оказалось < 5 , то следует присоединить к последнему.

Пример (С монеткой). Пусть $n = 4040$, $\#H = 2048$, $\#T = 1092$. Проверим $H_0 : \mathcal{P} = \text{Ber}(0.5)$ с $\alpha = 0.1$. Условия критерия выполняются, поэтому посчитаем

$$T = \frac{(2048 - 2020)^2}{2020} + \frac{(1092 - 2020)^2}{2020} = \frac{28^2 + 28^2}{2020} \approx 0.78 \sim \chi^2(1),$$

откуда

$$p\text{-value} = 1 - \text{cdf}_{\chi^2(1)}(0.78) \approx 0.38.$$

$0.38 > 0.1$, значит H_0 не отвергается.

Замечание. Прохождение критерия не достаточно. Так, альтернирующая (и явно не случайная) последовательность $\mathbf{x} = (0, 1, 0, 1, \dots)$ имеет $T = 0$.

4.4.3. Согласие с нормальным распределением

Для проверки гипотезы $H_0 : \mathcal{P}_\xi = N(a, \sigma^2)$ также можно воспользоваться статистикой критерия χ^2 для сложной гипотезы. В этом случае, нужно дискретизировать нормальное распределение, так, что

$$\mathcal{P}_0 = \begin{pmatrix} x_1^* & \dots & x_k^* \\ p_1(\hat{\boldsymbol{\theta}}) & \dots & p_k(\hat{\boldsymbol{\theta}}) \end{pmatrix}, \quad \hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}_{\text{MLE}}.$$

Тем не менее, нужно иметь в виду две теоретических неточности:

1. Построение \mathcal{P}_0 происходит случайно, в результате объединения элементов выборки после того, как она получена.
2. Оценка параметров $\hat{\theta}_{MLE}$ должна быть посчитана для \mathcal{P}_0 , а не для исходного (нормального) распределения — не $\bar{\mathbf{x}}, s^2$. Однако на практике на этот момент не обращают внимания.

Существует два возможных способа дискретизации:

1. Гистограмма: одинаковые интервалы, но разные вероятности.
2. Неравные интервалы с равными вероятностями.

```
N <- length(xs)
xs <- sort(xs)
probs <- pnorm(xs, mean=mean(xs), sd=sd(xs))
i <- 1; j <- i+1
while (N * (probs[j] - probs[i]) < 5) {
  j <- j+1
}
mean(xs[i:j]) # our  $x_1^*$ ;  $n_1 = j - i + 1$ 
```

Этот способ разбиения предпочтителен, потому что:

- Можно разбить максимально часто — так, чтобы $np_i = 5 \forall i$, следовательно и мощность будет максимальна.
- Он оказывается точнее первого на практике.
- Получается единственное p -value.

Замечание. Следует иметь в виду, что этот способ не годится для непрерывных, но плохо дискретизированных данных.

4.5. Критерий Колмогорова-Смирнова согласия с видом распределения

4.5.1. Произвольное абсолютно непрерывное распределение

$H_0 : \xi \sim \mathcal{P} = \mathcal{P}_0$.

Утверждение. Для проверки гипотезы согласия с видом произвольного абсолютно непрерывного распределения с известными параметрами используется асимптотический критерий Колмогорова-Смирнова со следующей статистикой:

$$D_n = \sup_{x \in \mathbf{x}} \left| \widehat{\text{cdf}}_n(x) - \text{cdf}_0(x) \right|,$$

где cdf_0 — функция распределения \mathcal{P}_0 нулевой гипотезы.

Альтернатива только одна: $H_1 : \xi \not\sim \mathcal{P}_0$; $\mathcal{A}_{\text{крит}} = (\text{qnt}_{\text{K-S}}(1 - \alpha), \infty)$.

Замечание. Критерий является не асимптотическим, но *точным*. Значит им пользоваться и при маленьких объемах выборки (мощность, при этом, останется низкой все-равно).

Замечание. $\sup_x \sqrt{n} \left| \widehat{\text{cdf}}_n(x) - \text{cdf}_0(x) \right| \xrightarrow{d} \mathcal{P}_{\text{K.S.}}$, где $\mathcal{P}_{\text{K.S.}}$ — распределение Колмогорова. Значит, при больших объемах выборки для такой статистики критерия можно пользоваться таблицами распределения Колмогорова.

Упражнение. Проверить гипотезу, что $\mathbf{x} = (0.1, 0.2, 0.4, 0.3, 0.1)$ есть выборка из $U[0, 1]$.

Решение. $D_n = 0.6$, p -value ≈ 0.05 (по таблицам или компьютером). Таким образом, при $\alpha > 0.05$ гипотеза отвергается, при $\alpha < 0.05$ — нет.

```
> ks.test(c(0.1,0.2,0.4,0.3,0.1), 'punif')
      One-sample Kolmogorov-Smirnov test
data:  c(0.1, 0.2, 0.4, 0.3, 0.1)
D = 0.6, p-value = 0.05465
```

┘

Замечание. Критерий Колмогорова-Смирнова консервативный — значение p -value завышено. Поэтому если гипотеза отвергается, то наверняка.

4.5.2. Нормальное распределение

Пусть $H_0 : \mathcal{P}_\xi \in \{N(a, \sigma^2)\}$. Как известно, критерий Колмогорова-Смирнова используется для непрерывных непараметрических распределений. Им можно воспользоваться и для данной H_0 , если вместо a, σ^2 подставить соответствующие оценки — в таком случае критерий будет консервативным. По аналогии с χ^2 хотелось бы сделать поправку на количество параметров — такая поправка осуществляется путем моделирования распределения тестовой статистики. Для $N(a, \sigma^2)$ и $\text{Exp}(\lambda)$ получаем распределение D_n , не зависящее от параметров (так что поправку можно делать вне зависимости от параметров; к примеру, $N(a, \sigma^2)$ можно привести к $N(0, 1)$ непрерывным преобразованием):

Критерий Бартлетта есть критерий Колмогорова-Смирнова для $H_0 : \mathcal{P}_\xi = \text{Exp}(\lambda)$.

Критерий Лиллиефорса¹⁰ для проверки $H_0 : \mathcal{P}_\xi = N(a, \sigma^2)$ считается статистикой D_n с $\text{cdf}_0(x) = \text{cdf}_{N(\bar{x}, s^2)}(x)$, сходящейся к распределению Лиллиефорса (Колмогорова-Смирнова с учетом подстановки оценок).

Критерий Шапиро-Уилка $T \approx \rho^2$, т.е. ведет себя примерно как квадрат коэффициента корреляции на normal probability plot.

Замечание. Распределения Лиллиефорса и Колмогорова-Смирнова были получены путем моделирования.

Пример. В R:

```
> require('nortest')
> lillie.test(rt(1000, df=10))
      Lilliefors (Kolmogorov-Smirnov) normality test
data:  rt(1000, df = 10)
D = 0.03042, p-value = 0.02953
```

4.6. Критерий типа ω^2

Определение. Статистика

$$Q = n \int_{\mathbb{R}} (\text{cdf}_n(x) - \text{cdf}_0(x))^2 w(x) d\text{cdf}_0(x),$$

где $w(x)$ — весовая функция.

Замечание. Статистика может быть проинтерпретирована как площадь разницы между соответствующими функциями распределения.

Cramer von Mises Q с $w \equiv 1$.

Anderson-Darling Q с

$$w(x) = \frac{1}{\text{cdf}_0(x)(1 - \text{cdf}_0(x))}.$$

Замечание. Весовая функция критерия Anderson-Darling присваивает большой вес значениям на хвостах распределения, поэтому сам критерий является мощным против разницы на хвостах, но и менее мощным при сдвиге.

Замечание. Все эти критерии точны.

Замечание. Распределение статистики в каждом случае не зависит от cdf_0 и все эти критерии состоятельные против любой альтернативы, поэтому не очень мощные.

4.7. Визуальное определение согласия с распределением

4.7.1. P-P plot

Определение. *P-P plot* есть график

$$\left\{ \left(\text{cdf}_0(x_i) + \frac{1}{2n}, \widehat{\text{cdf}}_n(x_i) \right) \right\}_{i=1}^n.$$

Пример. В R:

```
pp.plot <- function(xs, cdf.0=pnorm, n.knots=1000) {
  knots <- seq(min(xs), max(xs), length.out=n.knots)
  plot(cdf.0(knots), ecdf(xs)(knots))
  abline(0, 1)
}
```

4.7.2. Q-Q plot

Определение. *Q-Q plot* есть график

$$\left\{ \left(x_i, \text{cdf}_0^{-1} \left(\widehat{\text{cdf}}_n(x_i) + \frac{1}{2n} \right) \right) \right\}_{i=1}^n.$$

Определение. Частный случай Q-Q plot для $\text{cdf}_0^{-1} = \text{cdf}_{N(0,1)}^{-1}$ называется *normal probability plot*.

Пример. В R:

```
qq.plot <- function(xs, qf.0=qnorm, n.ppoints=1000) {
  qs <- ppoints(n.ppoints)
  plot(qf.0(qs), unname(quantile(xs, probs=qs)))
  abline(mean(xs), sd(xs))
}
```

Замечание. Если $\hat{\mathcal{P}}_n \rightarrow \mathcal{P}_\xi$, то оба графика будут стремиться к $y = x$. Референсной прямой normal probability plot будет $y = \widehat{D}\xi \cdot x + \widehat{E}\xi$.

Замечание. Больше о различии Q-Q и P-P plots, см. <http://v8doc.sas.com/sashtml/qc/chap8/sect9.htm>

Замечание. Различные интерпретации параметров распределения по Q-Q plot можно посмотреть в интерактивном приложении: <https://xiongge.shinyapps.io/QQplots/>

4.8. Гипотеза о равенстве распределений

$H_0 : \mathcal{P}_{\xi_1} = \mathcal{P}_{\xi_2}$.

Возможно рассматривать два случая:

Независимые выборки Две группы индивидов, на которых измеряется один и тот же признак. Формально: пусть $\zeta \in \{1, 2\}$ — номер группы, ξ — признак. Тогда $\xi_1 \sim \mathcal{P}_{\xi|\zeta=1}$, $\xi_2 \sim \mathcal{P}_{\xi|\zeta=2}$ и $\xi_1 \perp\!\!\!\perp \xi_2$. В этом случае выборка имеет вид

$$((x_1, x_2, \dots, x_{n_1}), (y_1, y_2, \dots, y_{n_2})).$$

Зависимые выборки Одна группа индивидов, на каждом из которых измеряются две характеристики (либо же «до» и «после»). В этом случае выборка имеет вид

$$((x_1, y_1), \dots, (x_n, y_n)).$$

Замечание. Для одной и той же гипотезы могут существовать разные критерии; их возможно сравнить по мощности, но только если они состоятельны против одной и той же альтернативы.

Замечание. Непараметрические критерии хороши тем, что основаны на рангах, значит устойчивы к аутлаерам; плохи тем, что не используют всю информацию о значении — только порядок, из-за чего обладают меньшей мощностью.

4.9. Равенство математических ожиданий для независимых выборок

4.9.1. Двухвыборочный t -критерий

$$H_0 : E\xi_1 = E\xi_2.$$

Определение. И для зависимых, и для независимых выборок используется *двухвыборочный t -критерий*

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{D(\bar{x} - \bar{y})}} \xrightarrow{\sim} N(0, 1).$$

Пусть выборка *независима*¹¹, $(x_1, \dots, x_{n_1}), (y_1, \dots, y_{n_2})$, $n = n_1 + n_2$. Значит $D(\bar{x} - \bar{y}) = D\bar{x} + D\bar{y}$.

Двухвыборочный t -критерий для независимых выборок с $\sigma_1^2 \neq \sigma_2^2$ (Welch t -test)

Предложение. Если дисперсия известна, $D(\bar{x} - \bar{y}) = D\bar{x} + D\bar{y} = \sigma_1^2/n_1 + \sigma_2^2/n_2$ и

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \xrightarrow{n_1, n_2 \rightarrow \infty} N(0, 1).$$

Если данные нормальные, то

$$t \sim N(0, 1).$$

Предложение. Если дисперсия неизвестна, $D(\widehat{\bar{x} - \bar{y}}) = s_1^2/n_1 + s_2^2/n_2$, откуда

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \xrightarrow{n_1, n_2 \rightarrow \infty} N(0, 1).$$

Замечание. Точное распределение неизвестно, примерно равно t с дробным числом степеней свободы (что вычисляется интерполяцией по соседним степеням). Всегда ожидается, что если данные нормальны, то распределение известно. Это противоречие носит название *проблемы Беренса-Фишера*¹².

¹¹Случай зависимой выборки рассматривается в другом параграфе.

¹²Behrens-Fisher problem.

Разбиение

$$H_1 : E\xi_1 \neq E\xi_2 \quad \mathcal{A}_{\text{крит}} = \mathbb{R} \setminus (z_{\alpha/2}, z_{1-\alpha/2})$$

$$H_1 : E\xi_1 > E\xi_2 \quad \mathcal{A}_{\text{крит}} = (z_{1-\alpha}, \infty)$$

$$H_1 : E\xi_1 < E\xi_2 \quad \mathcal{A}_{\text{крит}} = (-\infty, z_{\alpha})$$

Двухвыборочный t -критерий для независимых выборок с $\sigma_1^2 = \sigma_2^2$ (pooled t -test)

Предложение. Если дисперсия известна,

$$D(\bar{x} - \bar{y}) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} = \sigma^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right),$$

откуда

$$t = \frac{\bar{x} - \bar{y}}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \xrightarrow{n_1, n_2 \rightarrow \infty} N(0, 1).$$

Если данные нормальные, то

$$t \sim N(0, 1).$$

Предложение. Если дисперсия неизвестна,

$$t = \frac{\bar{x} - \bar{y}}{\tilde{s}_{1,2} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \xrightarrow{n_1, n_2 \rightarrow \infty} N(0, 1).$$

Если данные нормальные, то

$$t \sim t(n_1 + n_2 - 2).$$

Доказательство. Оценку дисперсии можно найти по объединенной и центрированной выборке (т.е. если H_0 верна, то $E\xi_1 = E\xi_2$ и можно думать как про одну выборку):

$$\begin{aligned} s_{1,2}^2 &= \frac{\overbrace{\sum_{i=1}^n (x_i - \bar{x})^2}^{\sim \chi^2(n_1-1)} + \overbrace{\sum_{i=1}^n (y_i - \bar{y})^2}^{\sim \chi^2(n_2-1)}}{n_1 + n_2} = \frac{n_1 \cdot s_1^2}{n_1 + n_2} + \frac{n_2 \cdot s_2^2}{n_1 + n_2} \\ \tilde{s}_{1,2}^2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (y_i - \bar{y})^2}{n_1 + n_2 - 2} = \frac{(n_1 - 1)\tilde{s}_1^2}{n_1 + n_2 - 2} + \frac{(n_2 - 1)\tilde{s}_2^2}{n_1 + n_2 - 2}, \end{aligned}$$

где в последнем случае оценка несмещенная и $E\tilde{s}_{1,2}^2 = \sigma^2$. □

Замечание. Этот вариант более точен, чем в случае $\sigma_1 \neq \sigma_2$.

Разбиение

$$H_1 : E\xi_1 \neq E\xi_2 \quad \mathcal{A}_{\text{крит}} = \mathbb{R} \setminus \left(\text{qnt}_{t(n_1+n_2-2)}(\alpha/2), \text{qnt}_{t(n_1+n_2-2)}(1 - \alpha/2) \right)$$

$$H_1 : E\xi_1 > E\xi_2 \quad \mathcal{A}_{\text{крит}} = (\text{qnt}_{t(n_1+n_2-2)}(1 - \alpha), \infty)$$

$$H_1 : E\xi_1 < E\xi_2 \quad \mathcal{A}_{\text{крит}} = (-\infty, \text{qnt}_{t(n_1+n_2-2)} \alpha)$$

Испытания Бернулли Пусть $\xi_i \sim \text{Ber}(p_i)$, $i \in \{1, 2\}$. Рассмотрим $H_0 : p_1 = p_2$ против $H_1 : p_1 \neq p_2$. Поскольку $E\xi_i = p_i$, применим двух-выборочный t -критерий. Объединим выборки и запишем:

$$D(\bar{x} - \bar{y}) = \frac{\hat{p}(1 - \hat{p})}{n_1 + n_2} \implies t = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1 - \hat{p})}} \xrightarrow{n_1, n_2 \rightarrow \infty} N(0, 1), \quad \hat{p} = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}.$$

Разбиение Аналогично с $\text{qnt}_{N(0,1)}$.

Определение. Выборка обладает *сбалансированным дизайном*, если $n_1 = n_2$.

Если дизайн сбалансирован, то

$$s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right) = \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2},$$

т.е. даже если дисперсии разные, результат одинаковый. Это остается справедливым даже при $n_1 \approx n_2$.

4.9.2. Непараметрический t -критерий

Можно использовать обычный t -критерий, но примененный к рангам.

Пусть, как и прежде, дана выборка (\mathbf{x}, \mathbf{y}) . Следующие два критерия — Wilcoxon и Mann-Whitney — проверяют гипотезу $H_0 : P(\xi_1 > \xi_2) = P(\xi_1 < \xi_2)$ или, альтернативно, что выборки получены из одной генеральной совокупности.

4.9.3. Критерии суммы рангов Wilcoxon

Следует сопоставить каждой выборке соответствующие её элементам ранги в *объединенной выборке*:

$$\begin{aligned} (x_1, \dots, x_{n_1}) &\mapsto (R_1, \dots, R_{n_1}) \\ (y_1, \dots, y_{n_2}) &\mapsto (T_1, \dots, T_{n_2}). \end{aligned}$$

Ясно, что если в целом элементы одной выборки окажутся больше другой, то нельзя будет говорить об их однородности. Определим

$$W_1 := \sum_{i=1}^{n_1} R_i, \quad W_2 := \sum_{i=1}^{n_2} T_i.$$

В качестве статистики можно было бы использовать либо W_1 , либо W_2 , однако, ни той, ни другой статистике невозможно априорно отдать предпочтение. Поэтому используется статистика

$$W := \max(W_1, W_2),$$

не имеющая аналитического выражения (но для которого посчитаны соответствующие таблицы).

Иногда в качестве статистики берут количество инверсий в объединенной выборке.

4.9.4. Критерий Mann-Whitney (U test)

Используется статистика

$$U := \max \left(n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - W_1, n_1 n_2 + \frac{n_2(n_2 + 1)}{2} - W_2 \right).$$

Замечание. $EU - n_1 n_2 / 2 \xrightarrow{n_1, n_2 \rightarrow \infty} 0$. Асимптотически,

$$\frac{U - EU}{\sqrt{DU}} \xrightarrow{n_1, n_2 \rightarrow \infty} N(0, 1),$$

но для малых объемов выборки можно посчитать и точные распределения.

Замечание. Критерий состоятельный против альтернативы

$$H_1 : P(\xi_1 > \xi_2) \neq P(\xi_1 < \xi_2).$$

Если формы распределений одинаковы, то эта альтернатива обозначает сдвиг. Для симметричных распределений это условие обозначает равенство медиан (а для нормального — математических ожиданий). Поэтому критерий устойчив к аутлаерам, хоть и за счет небольшой ($\approx 5\%$) потери мощности.

Замечание. Критерии Манна-Уитни и Вилкоксона *эквивалентны* — в том смысле, что выделяют один и тот же p -value. Тем не менее, проверяют они разные гипотезы ($E\xi$ не то же, что $\text{med } \xi$).

4.9.5. Критерий серий (runs)

Следует объединить выборку и в качестве статистики выбрать количество серий, т.е. подряд идущих элементов из одной выборки. Эта статистика имеет специально подобранное распределение.

4.9.6. Критерий равенства распределений

Рассматривается $H_0 : \mathcal{P}_{\xi_1} = \mathcal{P}_{\xi_2}$ против $H_1 : \mathcal{P}_{\xi_1} \neq \mathcal{P}_{\xi_2}$. Этот критерий мощный против неравенства распределений (а не отличия математического ожидания). В качестве статистики используется

$$D = \sup_x \left| \widehat{\text{cdf}}_{\xi_1}(x) - \widehat{\text{cdf}}_{\xi_2}(x) \right|.$$

Замечание. Все эти критерии подразумевают отсутствие повторяющихся наблюдений для избежания появления дробных рангов.

4.10. Равенство математических ожиданий для парных (зависимых) выборок

Выборка представлена набором пар $\{(x_i, y_i)\}_{i=1}^n$.

4.10.1. t -критерий

Пусть ξ_1, ξ_2 заданы на одном (Ω, \mathcal{F}, P) . Тогда гипотезу $H_0 : E\xi_1 = E\xi_2$ можно свести к $H_0 : E(\xi_1 - \xi_2) = E\eta = 0$ использовать не-парный t -тест.

Замечание (Мощность и зависимость). Сравним статистику для сбалансированного дизайна:

- Независимая выборка

$$t_{\text{indep}} = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{\sqrt{n}(\bar{x} - \bar{y})}{\sqrt{\sigma_1^2 + \sigma_2^2}}.$$

- Зависимая выборка:

$$\begin{aligned} D(\bar{x} - \bar{y}) &= D\bar{x} + D\bar{y} - 2 \text{cov}(\bar{x}, \bar{y}) = \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n} - 2\rho\sqrt{D\bar{x}}\sqrt{D\bar{y}} \\ &= \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n} - 2\rho\frac{\sigma_1}{\sqrt{n}}\frac{\sigma_2}{\sqrt{n}} = \frac{1}{n}(\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2), \end{aligned}$$

откуда

$$t_{\text{dep}} = \frac{\sqrt{n}(\bar{x} - \bar{y})}{\sqrt{\sigma_1^2 + \sigma_2^2 - 2\sigma_1\sigma_2\rho}}.$$

При $\rho > 0$, $t_{\text{dep}} > t_{\text{indep}}$. Значит, статистика чаще попадает в критическую область и критерий лучше находит различия (и мощность, следовательно, выше). Значит, тот же эксперимент на зависимых выборках мощнее.

Пример. Проверяют гипотезу, что белый свет лам влияет на решение задач.

- При тестировании на разных индивидах, должна быть уверенность, что они одинаковы по критичным параметрам (IQ, например).
- При тестировании на одинаковых индивидах следует составлять разные, но одинаковы по сложности задачи (второй раз одну и ту же задачу решать не займет много времени!). Мощность этого эксперимента будет выше.

4.10.2. Непараметрический тест знаков (Sign test)

$H_0 : P(\xi_1 < \xi_2) = P(\xi_1 > \xi_2)$. Используется статистика

$$W = \sum_{i=1}^n \psi_i, \quad \psi_i = \begin{cases} 1 & x_i > y_i \\ 0 & x_i < y_i. \end{cases}$$

Если при подсчете статистики $x_i = y_i$, эта пара игнорируется вместе с соответствующим уменьшением объема выборки.

Пусть после удаления всех пар, таких, что $x_i = y_i$, объем выборки стал равен m . Тогда $W \sim \text{Bin}(m, 0.5)$ и для построения разбиения можно пользоваться $\text{qnt}_{\text{Bin}(m, 0.5)}$.

Замечание. Критерий применим к порядковым признакам.

Замечание. Критерий очень устойчив к аутлаерам (но и очень низкомощен поэтому).

4.10.3. Непараметрический критерий (Paired Wilcoxon; Wilcoxon signed-rank test)

Увеличить мощность предыдущего критерия можно, учтя больше информации:

$$W := \sum_{i=1}^n R_i \psi_i, \quad R_i := \text{rk} |x_i - y_i|.$$

Для симметрии можно рассмотреть статистику

$$W = \sum_{i=1}^n R_i \text{sign}(x_i - y_i)$$

с идеальным значением 0. При верной H_0 , распределение W не имеет простого аналитического выражения (но может быть посчитана по таблицам), при этом $EW = 0$, $DW = n(n+1)(2n+1)/6$. Кроме того, $W \xrightarrow{d} N(0, DW)$, так что уже при $n \geq 10$ можно полагать, что $z = W/\sqrt{DW} \xrightarrow{d} N(0, 1)$ и строить разбиение соответственно.

Замечание. Критерий уже не применим к порядковым признакам.

4.11. Равенство дисперсии для двух распределений

$H_0 : D\xi_1 = D\xi_2$, $\xi_1 \perp\!\!\!\perp \xi_2$, $\xi_i \sim N(a, \sigma_i)^2$.

4.11.1. Критерий Фишера

$H_0 : \sigma_1^2 = \sigma_2^2$. Естественно использовать отношение s_1^2/s_2^2 с идеальным значением 1. Поделив на число степеней свободы, получим статистику

$$F := \frac{\tilde{s}_1^2}{\tilde{s}_2^2} \sim F(|\mathbf{x}| - 1, |\mathbf{y}| - 1).$$

Замечание. При отклонении от нормальности не становится асимптотическим.

4.11.2. Критерий Левена (Levene's test)

Так как $D\xi_i = E(\xi_i - E\xi_i)^2$, то критерий о равенстве дисперсий можно было бы свести к критерию о равенстве математических ожиданий; в этом случае применили бы t -критерий (подразумевающий разные дисперсии) к выборкам $\{(x_i - \bar{x})^2\}$ и $\{(y_i - \bar{y})^2\}$. Однако при возведении в квадрат распределение стало бы несимметричным и потребовался бы больший объем выборки. Кроме того, значительно бы усилились аутлаеры.

Вместо этого используют гипотезу $H_0 : E|\xi_1 - E\xi_1| = E|\xi_2 - E\xi_2|$ вместе с t -критерием, подразумевающим равенство дисперсий (для нормальных данных; иначе с разными).

4.11.3. Критерий Brown-Forsythe

Критерий Brown-Forsythe — это t -критерий для гипотезы $H_0 : E|\xi_1 - \text{med } \xi_1| = E|\xi_2 - \text{med } \xi_2|$.

Замечание. Устойчив к аутлаерам из-за использования $\text{med } \xi_i$.

5. Доверительное оценивание

5.1. Мотивация и определение

Для построенных оценок может понадобиться оценка точности. Так, даже состоятельная оценка может не быть в полном смысле «точной»: пусть $\theta_n^* \xrightarrow{P} \theta_0$; тогда

$$\hat{\theta}'_n = \begin{cases} c & n < N \gg 1 \\ \hat{\theta}_n & \text{иначе} \end{cases}$$

все-равно будет, конечно, состоятельной.

$D\hat{\theta}_n$ может быть не всегда просто вычислить и использовать.

Определение. $[c_1, c_2]$ — доверительный интервал для параметра θ_0 с уровнем доверия $\gamma \in [0, 1]$, если $\forall \theta_0$

$$P(\theta_0 \in [c_1, c_2]) = \gamma, \quad \text{где } c_1 = c_1(\mathbf{x}), c_2 = c_2(\mathbf{x}) \text{ — статистики.}$$

Замечание. Если выборка из дискретного распределения, то c_1, c_2 — тоже. Поэтому наперед заданную точность получить может не получиться; в таких случаях знак « $=$ » заменяют « \geq ». Аналогично с заменой на « $\xrightarrow[n \rightarrow \infty]{}$ ».

5.2. Доверительные интервалы для математического ожидания и дисперсии в нормальной модели

Предположение. Пусть $\xi \sim N(a, \sigma^2)$.

5.2.1. Доверительный интервал для a

- Пусть σ^2 известно. Свяжем a_0 с выборкой:

$$\gamma = P(c_1 < T < c_2) = P\left(c_1 < \sqrt{n} \frac{(\bar{x} - a_0)}{\sigma} < c_2\right) = P\left(a_0 \in \left(\bar{x} - \frac{\sigma c_2}{\sqrt{n}}, \bar{x} - \frac{\sigma c_1}{\sqrt{n}}\right)\right).$$

Решений уравнения $P(c_1 < \sqrt{n}(\bar{x} - a_0)/\sigma < c_2) = \Phi(c_2) - \Phi(c_1) = \gamma$ бесконечно много. Чем $[c_1, c_2]$ короче, тем лучше. Поскольку Φ симметрична и унимодальна,

$$\begin{aligned} c_1 &= -c_\gamma \\ c_2 &= c_\gamma, \end{aligned} \quad \text{где } c_\gamma = \text{cdf}_{N(0,1)}^{-1}\left(\gamma + \frac{1-\gamma}{2}\right) = x_{\frac{1+\gamma}{2}}.$$

Наконец,

$$P\left(a_0 \in \left(\bar{x} \pm \frac{\sigma}{\sqrt{n}} x_{\frac{1+\gamma}{2}}\right)\right) = \gamma.$$

- Пусть σ^2 неизвестно. По аналогии,

$$\gamma = P\left(c_1 < \frac{\sqrt{n-1}(\bar{x} - a_0)}{s} < c_2\right) = P\left(a_0 \in \left(\bar{x} \pm \frac{c_\gamma s}{\sqrt{n-1}}\right)\right), \quad c_\gamma = \text{cdf}_{t(n-1)}^{-1}\left(\frac{1+\gamma}{2}\right)$$

и

$$P\left(a_0 \in \left(\bar{x} \pm \frac{\tilde{s}}{\sqrt{n}} x_{\frac{1+\gamma}{2}}\right)\right) = \gamma.$$

Упражнение. Пусть $s^2 = 1.21$, $\bar{x} = 1.9$, $n = 36$. Построить 95% доверительный интервал для $E\xi$.

Решение.

$$c_\gamma = \text{qt}(0.975, 35) \approx 2.03 \implies \left(1.9 \pm \frac{2.03 \cdot \sqrt{1.21}}{\sqrt{35}}\right) = (1.52; 2.28).$$

┘

5.2.2. Доверительный интервал для σ^2

- Пусть a известно. Поскольку плотность χ^2 становится все более симметричной с ростом n , примем

$$c_1 = \text{cdf}_{\chi^2(n)}^{-1}\left(\frac{1-\gamma}{2}\right), \quad c_2 = \text{cdf}_{\chi^2(n)}^{-1}\left(\frac{1+\gamma}{2}\right).$$

Тогда

$$P\left(c_1 < \frac{ns_a^2}{\sigma_0^2} < c_2\right) = \gamma \iff P\left(\sigma_0^2 \in \left(\frac{ns_a^2}{x_{(1+\gamma)/2}}, \frac{ns_a^2}{x_{(1-\gamma)/2}}\right)\right) = \gamma.$$

- Пусть a неизвестно. Тогда аналогично

$$P\left(\sigma_0^2 \in \left(\frac{ns^2}{x_{(1+\gamma)/2}}, \frac{ns^2}{x_{(1-\gamma)/2}}\right)\right) = \gamma,$$

где $x_{(1\pm\gamma)/2} = \text{cdf}_{\chi^2(n-1)}^{-1}((1\pm\gamma)/2)$.

Определение. Случайная величина $g(x_1, \dots, x_n, \theta)$ называется *центральной статистикой параметра θ* , если

1. Её распределение («центральное распределение») не зависит от распределения θ .
2. G_n (функция распределения центрального распределения) непрерывна.
3. $\forall z_1, z_2$ и \mathcal{P}_θ -почти всюду

$$z_1 < g(x_1, \dots, x_n, \theta) < z_2$$

монотонно разрешимо относительно θ , т.е.

$$\exists f_1, f_n : f_1(x_1, \dots, x_n, \theta, z_1, z_2) < \theta < f_2(x_1, \dots, x_n, \theta, z_1, z_2).$$

Рассмотрим всегда разрешимое

$$\begin{aligned} \gamma &= G_n(z_2) - G_n(z_1) = P(z_1 < g(x_1, \dots, x_n, \theta) < z_2) \\ &= P(\underbrace{f_1(z_1, z_2, x_1, \dots, x_n)}_{c_1} < \theta < \underbrace{f_2(z_1, z_2, x_1, \dots, x_n)}_{c_2}). \end{aligned}$$

5.3. Асимптотический доверительный интервал для математического ожидания в модели с конечной дисперсией

Если модель неизвестна, но известно, что $D\xi < \infty$, можно построить доверительный интервал для $E\xi = a$. Пусть $\{x_i\}$ i.i.d., тогда

$$t = \frac{\sqrt{n}(\bar{x} - a)}{\sigma} \xrightarrow[n \rightarrow \infty]{} N(0, 1).$$

Если положить $\sigma := s$, то сходимость не испортится, потому что s^2 — состоятельная оценка σ^2 . Тогда

$$P\left(E\xi \in \left(\bar{x} \pm \frac{sc_\gamma}{\sqrt{n}}\right)\right) \xrightarrow[n \rightarrow \infty]{} \gamma, \quad c_\gamma = \text{cdf}_{t(n-1)}^{-1}\left(\frac{1+\gamma}{2}\right).$$

Альтернативно замену σ на s можно обосновать по теореме Slutsky.

Утверждение (Слутский). Если $\xi_n \xrightarrow{d} \xi$, $\eta_n \xrightarrow{P} c$, то $\xi_n + \eta_n \xrightarrow{d} \xi + c$ и $\xi_n \eta_n \xrightarrow{d} c\xi$.

Используя тот факт, что $s \xrightarrow{P} \sigma$, запишем

$$P\left(c_1 < \frac{\sqrt{n}(\bar{x} - a)}{\sigma} \frac{\sigma}{s} < c_2\right) \xrightarrow[n \rightarrow \infty]{} \Phi(c_2) - \Phi(c_1).$$

5.4. Асимптотический доверительный интервал для параметра на основе MLE

Если умеем находить $\hat{\theta}_{MLE}$, то по асимптотической нормальности,

$$\frac{\hat{\theta}_{MLE} - E\hat{\theta}_{MLE}}{\sqrt{D\hat{\theta}_{MLE}}} \xrightarrow{d} N(0, 1),$$

по асимптотической несмещенности,

$$\frac{\hat{\theta}_{MLE} - \theta}{\sqrt{D\hat{\theta}_{MLE}}} \xrightarrow{d} N(0, 1),$$

и, учитывая асимптотическую эффективность ($D\hat{\theta}_{MLE}I_n(\theta) \xrightarrow[n \rightarrow \infty]{} 1$), запишем статистику

$$T = (\hat{\theta}_{MLE} - \theta) \sqrt{I_n(\theta)} \xrightarrow{d} N(0, 1).$$

Чтобы по аналогии с предыдущим выразить θ в $P(c_1 < T < c_2) = P(|T| < c_\gamma) = \gamma$, необходимо выразить θ из $I_n(\theta)$. Для Pois и Ber это эквивалентно решению квадратного уравнения.

В общем случае, можно вместо θ в $I_n(\theta)$ подставить $\hat{\theta}_{MLE}$ (при $n \rightarrow \infty$ это не должно сильно испортить дело), откуда

$$P(|T| < c_\gamma) = \gamma \iff P\left(-c_\gamma < (\hat{\theta}_{MLE} - \theta) \sqrt{I_n(\theta)} < c_\gamma\right) = \gamma \iff P\left(\theta \in \left(\hat{\theta}_{MLE} \pm \frac{c_\gamma}{\sqrt{I_n(\theta)}}\right)\right) = \gamma,$$

где

$$T \xrightarrow{d} N(0, 1) \implies c_\gamma = \text{cdf}_{N(0,1)}^{-1}\left(\frac{1+\gamma}{2}\right).$$

Пример. $\xi \sim \text{Pois}(\lambda)$. По 2.5.1, $\hat{\lambda}_{MLE} = \bar{x}$, по 2.5.2 $I_n(\lambda) = n/\lambda = n/\bar{x}$ откуда

$$P\left(\lambda \in \left(\bar{x} \pm \text{cdf}_{N(0,1)}^{-1}\left(\frac{1+\gamma}{2}\right) \frac{\sqrt{\bar{x}}}{\sqrt{n}}\right)\right) = \gamma.$$

Пример. $\xi \sim \text{Ber}(p)$. $p = E\xi$. $\hat{p} = \bar{x}$, откуда

$$P\left(p \in \left(\hat{p} \pm c_\gamma \frac{\sqrt{\hat{p}(1-\hat{p})}}{\sqrt{n}}\right)\right) \xrightarrow[n \rightarrow \infty]{} \gamma.$$

Замечание. Этот доверительный интервал не очень хорош, потому что не принадлежит $[0, 1]$.

5.5. Доверительный интервал для проверки гипотезы о значении параметра

Зафиксируем $H_0 : \theta = \theta_0$ и $\gamma = 1 - \alpha$, где α играет роль уровня значимости. По определению доверительного интервала, $P(\theta \in [a_\gamma(\mathbf{x}), b_\gamma(\mathbf{x})]) = \gamma$. Тогда разбиением будет

$$\mathcal{A}_{\text{дов}} = [a_\gamma(\mathbf{x}), b_\gamma(\mathbf{x})], \quad \mathcal{A}_{\text{крит}} = \mathbb{R} \setminus \mathcal{A}_{\text{дов}},$$

причем

$$P(\theta_0 \notin [a_\gamma(\mathbf{x}), b_\gamma(\mathbf{x})]) = \alpha.$$

Иными словами, попадание в критическую область происходит с уровнем значимости α , что соответствует определению критерия.

5.5.1. Использование SE для построения доверительных интервалов

Пусть $\hat{\theta}$ — асимптотически нормальная оценка параметра θ . Чтобы сходимость оставалась верной, даже если подставляем оценку, нужна состоятельность. Тогда доверительный интервал для среднего с уровнем доверия γ имеет вид

$$E\hat{\theta} \pm c_\gamma \sqrt{D\hat{\theta}} \approx E\hat{\theta} \pm \text{cdf}_{N(0,1)}^{-1} \left(\frac{1+\gamma}{2} \right) \cdot \text{SE}.$$

6. Корреляционный и регрессионный анализы

Определение. Мера зависимости — это функционал $r : (\xi, \eta) \mapsto x \in [-1, 1]$ со свойствами:

1. $|r| \leq 1$.
2. $\xi \perp \eta \implies r(\xi, \eta) = 0$.
3. Если ξ и η «максимально зависимы», то $r(\xi, \eta) = 1$.

6.1. Вероятностная независимость

6.1.1. Визуальное определение независимости

- Поскольку при $p_\eta(y_0) \neq 0$

$$\xi \perp \eta \iff p_{\xi|\eta}(x | y_0) = \frac{p_{\xi,\eta}(x, y_0)}{p_\eta(y_0)} = p_\xi(x),$$

то срезы графика совместной плотности при фиксированном y_0 после нормировки $p_\eta(y_0)$ должны выглядеть одинаково для всех y_0 .

- Для выборки независимость можно попытаться определить по *таблицам сопряженности*: сгруппируем $\{(x_i, y_i)\}_{i=1}^n$ и сопоставим каждой уникальной паре абсолютную частоту ν_{ij} :

$$\begin{array}{cccc} & y_1^* & \cdots & y_s^* \\ x_1^* & \nu_{11} & \cdots & \nu_{1s} \\ \vdots & \vdots & \ddots & \vdots \\ x_k^* & \nu_{k1} & \cdots & \nu_{ks} \end{array}$$

Тогда признаки с большей чем случайной вероятностью будут независимы при пропорциональных строчках / столбцах. Более формально, признаки независимы, если

$$\frac{\nu_{ij}}{\sum_k \nu_{kj}} = \frac{\nu_{ij}}{\nu_{\cdot j}} = \hat{p}_{i|j} \propto \hat{p}_{i|\ell},$$

т.е. вероятности условного распределения не зависят от выбора строки.

Пример. Таблица сопряженности похожей на независимую выборки:

1	3	2
2	5	3
9	20	11

6.1.2. Критерий независимости χ^2

По определению, для двумерных дискретных распределений, независимость есть

$$\xi \perp\!\!\!\perp \eta \iff \underbrace{P(\xi = i, \eta = j)}_{p_{ij}} = \underbrace{P(\xi = i)}_{p_{i\cdot}} \underbrace{P(\eta = j)}_{p_{\cdot j}} = \underbrace{\sum_{k=1}^K P(\xi = i, \eta = k)}_{p_{i\cdot}} \cdot \underbrace{\sum_{s=1}^S P(\xi = s, \eta = j)}_{p_{\cdot j}}.$$

Проверим $H_0 : \xi \perp\!\!\!\perp \eta$.

Утверждение. ОМП оценкой будет $\hat{p}_{i\cdot} = \nu_{i\cdot}/n$ и $\hat{p}_{\cdot j} = \nu_{\cdot j}/n$.

Следовательно,

$$\xi \perp\!\!\!\perp \eta \iff \hat{p}_{ij} = \frac{\nu_{ij}}{n} = \hat{p}_{i\cdot} \hat{p}_{\cdot j} = \frac{\nu_{i\cdot}}{n} \cdot \frac{\nu_{\cdot j}}{n}.$$

Это равенство удается получить редко; важно определить, не является ли это нарушение случайным.

Запишем статистику

$$\chi^2 = \sum_{i=1}^K \sum_{j=1}^S \frac{(\nu_{ij} - n\hat{p}_{ij})^2}{n\hat{p}_{ij}} = \sum_{i=1}^K \sum_{j=1}^S \frac{(\nu_{ij} - \nu_{i\cdot}\nu_{\cdot j}/n)^2}{\nu_{i\cdot}\nu_{\cdot j}/n} \xrightarrow{d} \chi^2((k-1)(s-1))$$

Количество параметров таково, потому что если $\xi \parallel \eta$, то всего $ks - 1$ параметров (-1 потому что $\sum_{ij} p_{ij} = 1$); если $\xi \perp\!\!\!\perp \eta$, то $k + s - 2$ (-2 потому что $\sum_i p_{ij} = 1$ и $\sum_j p_{ij} = 1$). Значит $ks - 1 - k - s + 2 = (k-1)(s-1)$.

Пример. Дано S кубиков. Проверить гипотезу, что кубики одинаковы.

Решение. FIXME ┘

Замечание. На маленьких выборках ($n < 40$, $np_{ij} < 5$) возникают проблемы со сходимостью, потому что можно объединять только столбцы / строки и каждый раз терять сразу $S - 1$ ($K - 1$) степень свободы. В этих случаях используют критерием с перестановкой¹³ или, в случае таблиц сопряженности 2×2 , точным критерием Фишера.

Замечание. Критерий верен для количественных, порядковых и качественных признаков, потому что нигде не участвуют значения из выборки.

Замечание. Критерий асимптотический, поэтому $\alpha_1 \rightarrow \alpha$.

Замечание. Критерий не удовлетворяет 1-му пункту определения меры зависимости ($\chi^2 \notin [-1, 1]$). Это обычно исправляют так: рассматривают *среднеквадратичную сопряженность*

$$r^2 := \frac{\chi^2}{n}$$

или коэффициент сопряженности Пирсона

$$p^2 := \frac{\chi^2}{\chi^2 + n}$$

(тогда 1 никогда не достигается). Могли бы работать с $1 - p\text{-value}$, но так почему-то никогда не делают.

¹³[https://en.wikipedia.org/wiki/Resampling_\(statistics\)#Permutation_tests](https://en.wikipedia.org/wiki/Resampling_(statistics)#Permutation_tests)

6.2. Линейная / полиномиальная зависимость

Пусть теперь ξ, η — количественные признаки.

Определение. Определим

$$\phi(x) := E\{\eta \mid \xi = x\}.$$

Тогда назовем зависимость *линейной*, если $\phi(x)$ — линейная функция, *квадратичной* — если квадратичная и т.д.

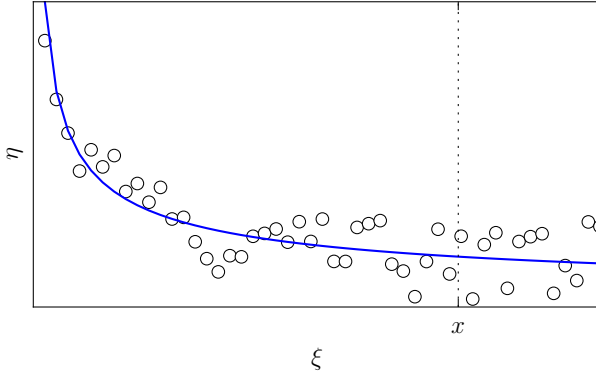


Рис. 3: Нелинейная зависимость

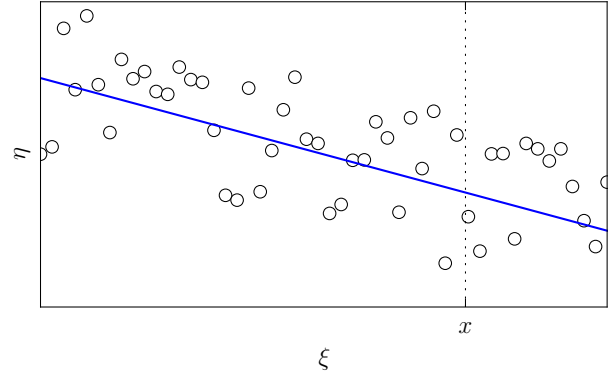


Рис. 4: Линейная зависимость

Определение. Мера *линейной* зависимости между случайными величинами ξ и η есть *коэффициент корреляции Пирсона*

$$\rho = \frac{\text{cov}(\xi, \eta)}{\sqrt{D\xi}\sqrt{D\eta}}.$$

Замечание. Про ρ можно думать как про \cos между векторами в соответствующем пространстве.

Замечание (Важное).

$$\begin{aligned} \xi \perp\!\!\!\perp \eta &\implies \rho = 0 \\ \xi, \eta \sim N(a, \sigma^2), \xi \perp\!\!\!\perp \eta &\iff \rho = 0. \end{aligned}$$

Предложение. Для линейно зависимых данных, конечно, $\rho = 1$.

Доказательство. Пусть $\eta = a + b\xi$; тогда

$$\begin{aligned} \rho(\xi, \eta) &= \frac{\text{cov}(\xi, a + b\xi)}{\sqrt{D\xi}\sqrt{D(a + b\xi)}} = \frac{E\xi(a + b\xi) - E\xi E(a + b\xi)}{\sqrt{D\xi}\sqrt{Db\xi}} = \frac{E\xi a + bE\xi^2 - E\xi Ea - E\xi bE\xi}{b\sqrt{D\xi}\sqrt{D\xi}} = \\ &= \frac{aE\xi + bE\xi^2 - aE\xi - b(E\xi)^2}{bD\xi} = \frac{b(E\xi^2 - (E\xi)^2)}{bD\xi} = 1. \end{aligned}$$

□

О соотношении ρ и коэффициента линейной регрессии По (6.5.1), если линейная регрессия уравнением $y = kx + b$, то

$$k = \rho \frac{\sigma_\eta}{\sigma_\xi}.$$

В общем случае, по виду прямой линейной регрессии ничего нельзя сказать о зависимости между случайными величинами. Так, если $\eta = a + b\xi$ есть линейная функция от ξ , то, по предыдущему, $\rho = 1$ и

$$k = 1 \cdot \frac{\sqrt{D(a + b\xi)}}{\sqrt{D\xi}} = b$$

и прямая может иметь произвольный, в зависимости от b , наклон.

Замечание. В то же время, поскольку для

$$\begin{pmatrix} \xi \\ \eta \end{pmatrix} \sim N(\boldsymbol{\mu}, \Sigma), \quad \Sigma = \begin{pmatrix} \sigma_\xi^2 & \text{cov}(\xi, \eta) \\ \text{cov}(\xi, \eta) & \sigma_\eta^2 \end{pmatrix}$$

справедливо, что

$$k = \rho \frac{\sigma_\eta}{\sigma_\xi} = \frac{\text{cov}(\xi, \eta)}{\sigma_\xi \sigma_\eta} \cdot \frac{\sigma_\eta}{\sigma_\xi} = \frac{1}{\sigma_\xi^2} \text{cov}(\xi, \eta),$$

то $k = 0 \iff \text{cov}(\xi, \eta) = 0$, а для стандартно нормальных данных $k = \rho = \text{cov}(\xi, \eta)$.

Значимость коэффициента корреляции

Определение. Коэффициент корреляции *значим*, если отвергается $H_0 : \rho = 0$.

Статистика для проверки значимости при $(\xi, \eta)^T \sim N(\boldsymbol{\mu}, \Sigma)$ (иначе статистика асимптотическая)

$$T = \frac{\sqrt{n-2} \hat{\rho}_n}{\sqrt{1 - \hat{\rho}_n^2}} \sim t(n-2).$$

Идеальное значение — 0, два хвоста.

6.3. Метод наименьших квадратов (Ordinary Least Squares)

Пусть $\eta, \xi \in L^2(\mathcal{F}, \mathbf{P})$ пространству \mathcal{F} -измеримых по мере \mathbf{P} функций с конечным вторым моментом и скалярным произведением $(\eta, \xi) = \mathbf{E} \eta \xi$, причем $\hat{\eta} \in K = \{\phi(\xi)\} = \{\hat{\eta} : \sigma(\phi(\xi))\text{-измерима}\}$. По свойству УМО(2.3.2), вектор

$$\hat{\eta}^* = \mathbf{E} \{\eta \mid \phi(\xi)\}$$

будет ортогональной проекцией η на K , т.е. $(\eta - \hat{\eta}^*, \hat{\eta}) = 0 \forall \hat{\eta} \in K$. Значит, он минимизирует квадрат нормы расстояния от η до K :

$$\hat{\eta}^* = \underset{\hat{\eta} \in K}{\operatorname{argmin}} \|\eta - \hat{\eta}\|^2 = \underset{\hat{\eta} \in K}{\operatorname{argmin}} \mathbf{E} (\eta - \hat{\eta})^2 = \mathbf{E} \{\eta \mid \phi(\hat{\eta})\}.$$

$\hat{\eta}^*$ называется *наилучшим среднеквадратичным приближением в классе K* .

6.4. Корреляционное отношение

Если $K = \mathcal{L} = \{a\xi + b\}$ — линейное пространство, то теорема Пифагора принимает вид

$$D\eta = \mathbf{E} (\eta - \mathbf{E}\eta)^2 = \underbrace{\mathbf{E} (\hat{\eta}^* - \mathbf{E}\eta)^2}_{\text{объяснённая доля аппроксимации}} + \underbrace{\mathbf{E} (\eta - \hat{\eta}^*)^2}_{\text{ошибка аппроксимации}}.$$

Откуда можно записать меру аппроксимации

$$\frac{\mathbf{E} (\hat{\eta}^* - \mathbf{E}\eta)^2}{D\eta} = 1 - \frac{\mathbf{E} (\eta - \hat{\eta}^*)^2}{D\eta} = 1 - \frac{\min_{\hat{\eta} \in \mathcal{L}} \mathbf{E} (\eta - \hat{\eta})^2}{D\eta}.$$

Определение. Полученная величина называется коэффициентом корреляции ρ^2 :

$$\rho^2 := 1 - \frac{\min_{\hat{\eta} \in \mathcal{L}} \mathbf{E} (\eta - \hat{\eta})^2}{D\eta}.$$

ρ — коэффициент корреляции Пирсона.

Определение. Множественный коэффициент корреляции есть полученная величина для МНК с $K = \mathcal{M} = \left\{ \sum_{i=1}^k b_i \xi_i + b_0 \right\}$.

$$R^2(\eta, \xi_1, \dots, \xi_k) := 1 - \frac{\min_{\hat{\eta} \in \mathcal{M}} \mathbf{E} (\eta - \hat{\eta})^2}{D\eta}.$$

Замечание. $R^2 \geq \rho^2$; если же $R^2 = \rho^2$, то ξ_1, \dots, ξ_k все зависимы.

Определение. В общем случае, если $K = \{\phi(\xi) \text{ измеримые}\}$, то полученная величина называется *корреляционным отношением*:

$$r_{\eta|\xi}^2 := 1 - \frac{\min_{\hat{\eta} \in K} \mathbb{E}(\eta - \hat{\eta})^2}{D\eta} = \frac{D\mathbb{E}(\eta | \xi)}{D\eta}.$$

Свойства корреляционного отношения

1. $r_{\eta|\xi}^2 \in [0, 1]$.
2. $\eta \perp \xi \implies r_{\eta|\xi}^2 = 0$.
3. $\eta = \phi(\xi) \iff r_{\eta|\xi}^2 = 1$.
4. Вообще говоря, $r_{\eta|\xi}^2 \neq r_{\xi|\eta}^2$. К примеру, для любой не монотонной функции (так, чтобы не существовала обратная).
5. $r_{\eta|\xi}^2 \geq \rho^2(\eta, \xi)$ (потому что минимум по всем функциям меньше, чем лишь по линейным, значит $1 - \min$ больше).
6. $(\xi, \eta)^T \sim N(\mu, \Sigma) \implies r_{\eta|\xi}^2 = \rho^2(\eta, \xi)$.

Выборочное корреляционное отношение

$$\hat{r}_{\eta|\xi}^2 = \hat{r}_{y|x}^2 = \frac{s_{y|x}^2}{s_y^2}.$$

6.5. Регрессия

Определение. Регрессией η по ξ называется $\mathbb{E}\{\eta | \xi\}$.

Замечание. Таким образом осуществляется предсказание η по ξ с минимальной среднеквадратичной ошибкой.

Определение. Функция регрессии есть $f(x) = \mathbb{E}\{\eta | \xi = x\}$.

Замечание. f находится по МНК для $K = \{\psi(\xi) : \psi \text{—измеримая}\}$.

Виды регрессий

- Нелинейными и линейными ($K = \{a\xi + b\}$);
- Парными (предсказывая величину по одной случайной величине) и множественными (по многим).

6.5.1. Парная линейная регрессия

Определение. Пусть $\xi, \eta \in L^2$. Парной линейной регрессией η по ξ называется наилучшее среднеквадратичное приближение $h_{\beta_1^*, \beta_2^*}(\xi) = \beta_1^* \xi + \beta_2^*$ в классе линейных по ξ функций $K = \mathcal{L} = \{\beta_1 \xi + \beta_2\}$. Иными словами,

$$h_{\beta_1^*, \beta_2^*}(\xi) = \operatorname{argmin}_{\beta_1, \beta_2} \|\eta - h_{\beta_1, \beta_2}(\xi)\|^2 = \mathbb{E}\{\eta | h_{\beta_1, \beta_2}(\xi)\} = \operatorname{argmin}_{\beta_1, \beta_2} \underbrace{\mathbb{E}(\eta - (\beta_1 \xi + \beta_2))^2}_{\phi(\beta_1, \beta_2)} = \beta_1^* \xi + \beta_2^*.$$

Замечание. Найти минимум ϕ можно, как обычно, решив систему $\partial\phi/\partial\beta_i = 0$ ¹⁴.

Утверждение. β_1^*, β_2^* таковы, что

$$\frac{h(\xi) - E\eta}{\sqrt{D\eta}} = \rho \frac{\xi - E\xi}{\sqrt{D\xi}}.$$

Это уравнение задает линию регрессии. Иными словами,

$$h(\xi) = \underbrace{\rho \frac{\sqrt{D\eta}}{\sqrt{D\xi}}}_{\beta_1^*} \xi + \underbrace{E\eta - \rho \frac{\sqrt{D\eta}}{\sqrt{D\xi}} E\xi}_{\beta_2^*}.$$

Отсюда можно получить соотношение между коэффициентом линейной регрессии $\beta_1^* = k$ (наклоном регрессионной прямой) и коэффициентом корреляции:

$$k = \rho \frac{\sigma_\eta}{\sigma_\xi}.$$

Замечание. Подстановкой проверяется, что

$$\phi(\beta_1^*, \beta_2^*) = \min_{\hat{\eta} \in K} E(\eta - \hat{\eta})^2 = D\eta(1 - \rho^2),$$

откуда можно найти уже известное выражение для коэффициента корреляции Пирсона

$$\rho^2(\eta, \xi) = 1 - \frac{\phi(\beta_1^*, \beta_2^*)}{D\eta} = 1 - \frac{\min_{\hat{\eta} \in H} E(\eta - \hat{\eta})^2}{D\eta}, \quad \hat{\eta} := h(\xi).$$

Определение. Линейная регрессия *значима*, если $\beta_1^* \neq 0 \implies \rho \neq 0$. Значимость регрессии эквивалентна значимости предсказания по ней.

Определение. Величина *sum of squares residual* есть

$$SSR = n \cdot \phi(\beta_1^*, \beta_2^*) = \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad \hat{y}_i = h_{\beta_1^*, \beta_2^*}(x_i).$$

6.5.2. Модель линейной регрессии

Можно описать выборку как

$$y_i = \beta_1 x_i + \beta_2 + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2), \quad \epsilon_i \perp \epsilon_j.$$

σ^2 — мешающий параметр, который можно оценить через SSR/n . Но если $\epsilon_i \sim N(0, \sigma^2)$, то

$$\hat{\sigma}^2 = \frac{SSR}{n-2}$$

есть несмещенная оценка σ^2 . Значит,

$$\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sigma^2} \sim \chi^2(n-2).$$

Замечание. МНК минимизирует разницу $y_i - \hat{y}_i$, что на графике соответствует вертикальным отрезкам, соединяющим y_i и $\hat{y}_i = h(x_i)$. Это не то же, что минимизация перпендикуляров от y_i на $h(x)$ — техники метода анализа главных компонент («РСА»).

Замечание. Существует три базовых модели, в которых функция регрессии линейная:

- $\eta = \beta_1 \xi + \beta_2 + \epsilon$, $\epsilon \perp \xi$, $E\epsilon = 0$.
- $(\xi, \eta)^T \sim N(\mu, \sigma^2)$.
- ξ принимает всего два значения (возможно, как качественный признак).

¹⁴См. https://en.wikipedia.org/wiki/Simple_linear_regression

6.5.3. Доверительные интервалы для β_1 и β_2

Как обычно, помимо точечной оценки $\hat{\beta}_1$ и $\hat{\beta}_2$, интересуемся диапазоном значений, которые может принимать оценка с заданной вероятностью. Примем предположение о несмещенности оценки, т.е. $E\hat{\beta}_i = \beta_i$. Поскольку в модели $y_i = \beta_1 x_i + \beta_2 + \epsilon_i$ ошибка $\epsilon_i \sim N(0, \sigma^2)$ есть случайная величина, оценки $\hat{\beta}_i$ — тоже становятся случайными величинами: $\hat{\beta}_i \sim N(\beta_i, D\hat{\beta}_i)$. В курсе регрессионного анализа доказывается¹⁵, что

$$D\hat{\beta}_1 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad D\hat{\beta}_2 = \frac{\sigma^2}{n}.$$

Кроме того,

$$SE(\hat{\beta}_1) = \sqrt{D\hat{\beta}_1} = \frac{\hat{\sigma}}{\sqrt{ns_x}} = \frac{\sqrt{\frac{SSR}{n-2}}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}, \quad SE(\hat{\beta}_2) = SE(\hat{\beta}_1) \cdot s_x$$

Предложение. *Статистика*

$$t = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t(n-2).$$

Доказательство. Известно,

$$t \sim t(m) \iff t = \frac{\xi}{\sqrt{\eta/m}}, \quad \xi \sim N(0, 1), \quad \eta \sim \chi^2(m).$$

Ясно, что

$$\frac{\hat{\beta}_1 - \beta_1}{\frac{\sigma}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim N(0, 1), \quad \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sigma^2} \sim \chi^2(n-2).$$

Тогда

$$\frac{\left(\frac{\hat{\beta}_1 - \beta_1}{\left(\frac{\sigma}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \right)} \right)}{\left(\frac{\left(\frac{\sqrt{\sum_{i=1}^n (y_i - \hat{y}_i)^2}}{\sigma} \right)}{\sqrt{n-2}} \right)} = \frac{(\hat{\beta}_1 - \beta_1) \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}{\frac{\sigma}{\sqrt{\sum_{i=1}^n \hat{\epsilon}_i^2}}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t(n-2).$$

□

Используя статистику t , введем доверительные интервалы с $c_\gamma = \text{cdf}_{t(n-2)}^{-1}((1+\gamma)/2)$:

$$t \in (-c_\gamma, c_\gamma) \iff \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \in (-c_\gamma, c_\gamma) \iff \beta_1 \in \left(\hat{\beta}_1 - c_\gamma SE(\hat{\beta}_1), \hat{\beta}_1 + c_\gamma SE(\hat{\beta}_1) \right).$$

Аналогично, для β_2 :

$$\beta_2 \in \left(\hat{\beta}_2 - c_\gamma SE(\hat{\beta}_2), \hat{\beta}_2 + c_\gamma SE(\hat{\beta}_2) \right).$$

¹⁵См. https://en.wikipedia.org/wiki/Proofs_involving_ordinary_least_squares

Замечание. На картинке доверительные интервалы изображаются в виде «рукавов» вокруг графика линейной регрессии — т.е. область всевозможных положений прямой при варьировании β_1, β_2 в заданных интервалах.

Пример. Линейная регрессия как предсказательная модель может быть использована неправильно в следующих случаях:

- неправильная модель;
- применение к неоднородным данным (аутлаер или неоднородность);
- хотим построить предсказание в точке, далекой от данных (проблема — большая ошибка);
- не знаем какая модель там, где данных нет.

6.6. Частная корреляция

Определение. Частная корреляция случайных величин η_1, η_2 относительно $\{\xi_1, \dots, \xi_k\}$ есть

$$\rho(\eta_1, \eta_2 \mid \{\xi_1, \dots, \xi_k\}) := \rho(\eta_1 - \hat{\eta}_1^*, \eta_2 - \hat{\eta}_2^*), \quad \text{где } \hat{\eta}_i^* = \underset{\hat{\eta}_i \in \{\sum_{i=1}^k b_i \xi_i + b_0\}}{\operatorname{argmin}} \mathbb{E}(\eta_i - \hat{\eta}_i)^2.$$

Если регрессия линейна, то

$$\rho(\eta_1, \eta_2 \mid \xi_1, \dots, \xi_k) = \rho(\eta_1 - \mathbb{E}\{\eta_1 \mid \xi_1, \dots, \xi_k\}, \eta_2 - \mathbb{E}\{\eta_2 \mid \xi_1, \dots, \xi_k\}).$$

Замечание (Важное). Пусть в эксперименте подсчитан ненулевой ρ . Это может означать, что один из факторов является причиной, а другой следствием; чтобы установить, что есть что, проводят эксперимент и смотрят, какой фактор в реальности влияет на какой. Это может также означать, что влияет сторонний фактор. Чтобы его исключить, считают частную корреляцию.

Пример. Возможна ситуация, когда $\rho(\eta_1, \eta_2) \neq 0$, но $\rho(\eta_1, \eta_2 \mid \xi) = 0$. Частная корреляция есть, по сути, корреляция на центрированных данных.

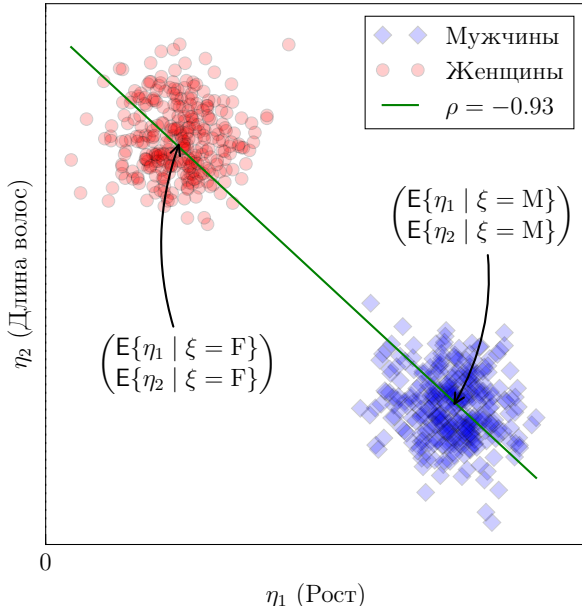


Рис. 5: Исходные данные (бимодальность)

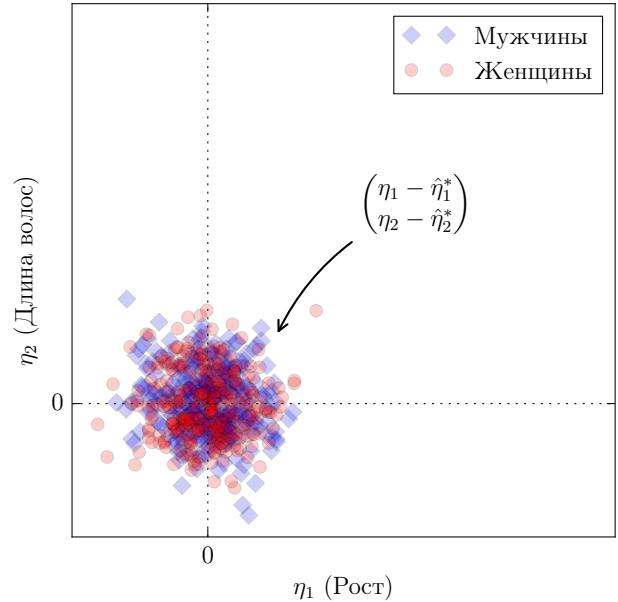


Рис. 6: Центрированные данные

Пример. Возможна и ситуация как на (7), где определено $\rho(\eta_1, \eta_2) > 0$, но $\rho(\eta_1, \eta_2 \mid \xi) < 0$.

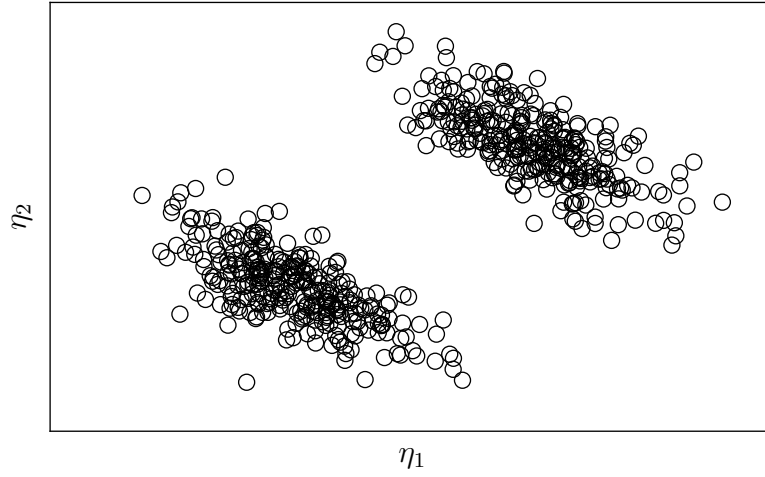


Рис. 7: $\rho(\eta_1, \eta_2) > 0$, но $\rho(\eta_1, \eta_2 \mid \xi) < 0$

Замечание. По аналогии с предыдущим примером, если $|\text{im } \xi| \rightarrow \infty$, то графики (η_1, η_2) при фиксированном ξ образуют эллипсоид (в этом случае с положительной корреляцией).

6.7. Зависимость между порядковыми признаками

Пусть на выборке

$$\begin{pmatrix} \xi \\ \eta \end{pmatrix} \sim \begin{pmatrix} x_1 \\ y_1 \end{pmatrix}, \dots, \begin{pmatrix} x_n \\ y_n \end{pmatrix}$$

задан только порядок. Тогда можем считать только эмпирическую функцию распределения.

Следующие коэффициенты основаны на рангах. Ранговые характеристики хорошо работают на выборках *без повторений* (чтобы, к примеру, не возникало дробных рангов).

6.7.1. Ранговый коэффициент Спирмана

Определение. Ранговый коэффициент Спирмана есть

$$\rho_S = \rho(\text{cdf}_\xi(\xi), \text{cdf}_\eta(\eta)).$$

Замечание. $\text{cdf}_\xi(\xi) \sim U(0, 1)$, потому что $P(\text{cdf}_\xi(\xi) < x) = P(\xi < \text{cdf}_\xi^{-1}(x)) = \text{cdf}_\xi(\text{cdf}_\xi^{-1}(x)) = x$.

Определение. Ранг элемента из выборки есть его порядковый номер в упорядоченной выборке:

$$\text{rk } x_{(i)} = i.$$

Обозначение. $\text{rk } x_{(i)} =: R_i$, $\text{rk } y_{(i)} =: T_i$.

Можем ввести эмпирическое распределение

$$\text{cdf}_{\xi_n}(x_i) = \frac{\text{rk } x_i}{n}, \quad \text{cdf}_{\eta_n}(y_i) = \frac{\text{rk } y_i}{n} = \frac{T_i}{n}.$$

Тогда будет справедливо следующее

Определение. Выборочный коэффициент Спирмана определяется как выборочный коэффициент корреляции Пирсона $\hat{\rho}$, но с заменой значений на ранги:

$$\hat{\rho}_S \begin{pmatrix} \xi_n \\ \eta_n \end{pmatrix} = \rho \begin{pmatrix} R_n \\ T_n \end{pmatrix} = \frac{1/n \cdot \sum_{i=1}^n R_i T_i - \bar{R} \bar{T}}{\sqrt{1/n \cdot \sum_{i=1}^n (R_i - \bar{R})^2} \sqrt{1/n \cdot \sum_{i=1}^n (T_i - \bar{T})^2}}.$$

Если нет повторяющихся наблюдений, то знаменатель будет одним и тем же у всех выборок объема n , значит его можно посчитать заранее. В этом (и только этом) случае, справедлива более простая формула:

$$\hat{\rho}_S = 1 - \frac{6 \sum_{i=1}^n (R_i - T_i)^2}{n^3 - n}.$$

Замечание. Из последней формулы хорошо видно, что если x_i, y_i все идут в одном порядке, то $R_i - T_i = 0 \forall i$ и $\hat{\rho}_S = 1$.

Замечание. ρ_S для количественных признаков есть мера монотонной зависимости:

$$\rho_S = 1 \iff (x_i > x_{i+1} \implies y_i > y_{i+1} \forall i)$$

(даже если зависимость нелинейная и $\rho \neq 1$). Иными словами, $\rho_S > 0$, если y имеет тенденцию к возрастанию с возрастанием x (и $\rho_S < 0$ иначе). Чем большее ρ_S , тем более явно выражена зависимость y от x в виде некоторой монотонной функции.

Согласованность ρ и ρ_S ρ_S не согласована с ρ в том же смысле, что ρ и $r_{\xi|\eta}$.

Утверждение. Если данные нормальные то справедлива формула

$$\rho = 2 \sin \left(\frac{\pi}{6} \rho_S \right).$$

Значит, можем сравнить критерии между собой.

- С точностью до погрешности, по значению, $\hat{\rho}$ и $\hat{\rho}_S$ — это одно и то же (см. 8)

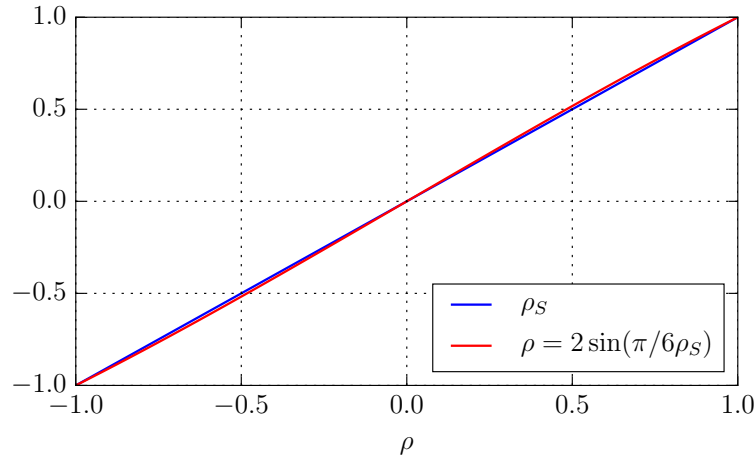


Рис. 8: $\hat{\rho} \approx \hat{\rho}_S$

- Обычный критерий оценки — выборочную дисперсию — посчитать сложно. Тем не менее, можем заметить, что $\hat{\rho}_S$ более устойчив к аутлаерам (см. 9). Всегда можно добавить аутлаер такой, что $\hat{\rho} = 0$; $\hat{\rho}_S$ же поменяется не сильно. Поэтому для нормальных данных, ρ_S — это оценка, что нет аутлаеров.

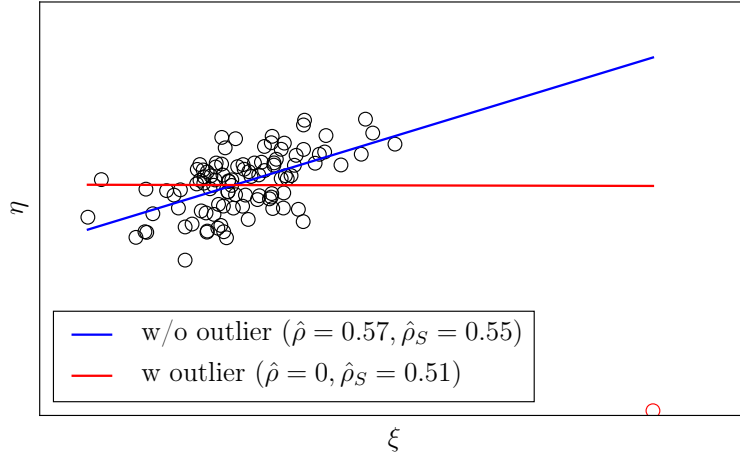


Рис. 9: $\hat{\rho}$ до и после добавления аутлаера

- Монотонным преобразованием можем всегда сделать так, чтобы ρ изменился (например, возведя в квадрат); при монотонном преобразовании, однако, не меняется ρ_S (см. 10). Значит, чтобы узнать ρ исходных (нормальных) данных, можно не выполнять обратного преобразования, а сразу посчитать ρ_S .

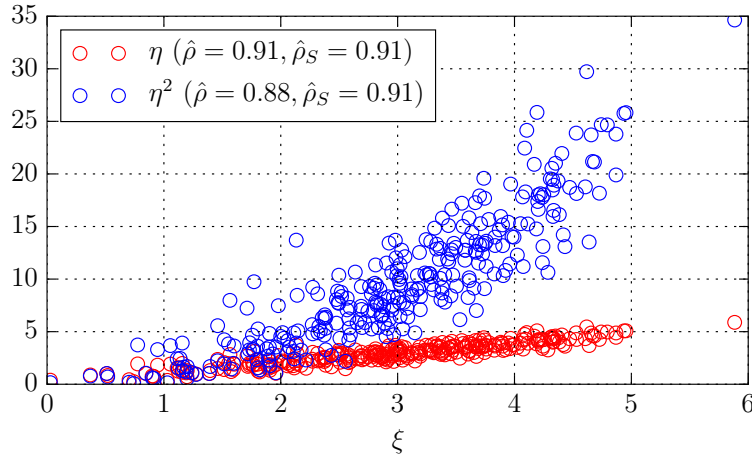


Рис. 10: Монотонное преобразование нормальных данных

6.7.2. Ранговый коэффициент Кэндалла $\tau(\xi, \eta)$

Определение. Пусть $(\xi_1, \eta_1)^T \perp (\xi_2, \eta_2)^T \sim \mathcal{P}_{\xi, \eta} \sim (\xi, \eta)^T$; тогда *ранговым коэффициентом Кэндалла* называется

$$\tau(\xi, \eta) = \rho(\text{sign}(\xi_2 - \xi_1), \text{sign}(\eta_2 - \eta_1)) = P((\xi_2 - \xi_1)(\eta_2 - \eta_1) > 0) - P((\xi_2 - \xi_1)(\eta_2 - \eta_1) < 0).$$

На выборочном языке, пусть дана выборка $(x_1, y_1), \dots, (x_n, y_n)$; тогда

$$\tau = \frac{\#(\text{одинаково упорядоченных пар}) - \#(\text{по-разному упорядоченных пар})}{\#(\text{комбинаций пар})},$$

где пара $(x_i, y_i), (x_j, y_j)$ считается одинаково упорядоченной, если $\text{sign}(x_i - x_j) = \text{sign}(y_i - y_j)$, а $\#(\text{комбинаций пар}) = C_n^2 = n(n-1)/2$.

Утверждение. Если $(\xi, \eta)^T \sim N(\boldsymbol{\mu}, \Sigma)$, то справедлива формула

$$\rho = \sin\left(\frac{\pi}{2}\tau\right).$$

Из утверждения следует, что τ все время меньше ρ и ρ_S (по модулю).

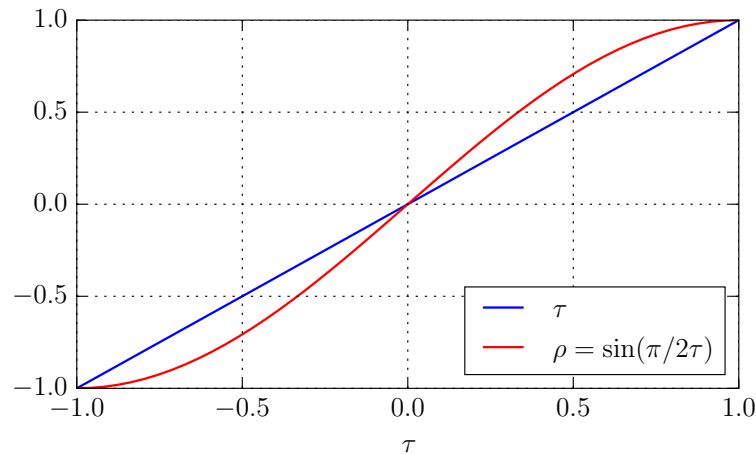


Рис. 11: ρ и τ

Пример (Проверка ряда на тренд). Пусть ξ — номера точек, а η — значения ряда. Тогда H_0 : $\tau_0 = 0$ и если H_0 отвергается, то тренд присутствует.

6.8. Корреляционные матрицы

Если признаков много, то их наглядно характеризуют корреляционные матрицы. Улучшить наглядность можно переупорядочив признаки так, чтобы на диагонали матрицы стояли блоки корреляций признаков из «корреляционных плеяд».

Определение. Пусть ρ_0 ; корреляционная плеяда есть множество признаков, таких, что их попарная корреляция больше ρ_0 .

Можно выделить и несколько уровней ρ_i : $\rho_0 < \rho_1 < \dots$. Тогда сначала следует составить плеяду по ρ_0 , затем внутри полученного по ρ_1 и т.д.

А. Другие полезные распределения случайных величин

А.1. Пуассона

А.2. Логнормальное

В. Свойства условного математического ожидания

1. $E\{a\xi + b\theta \mid \eta\} = aE\{\xi \mid \eta\} + bE\{\theta \mid \eta\}.$

2. $E E\{\eta \mid \xi\} = E\eta.$

3. $\xi \perp \eta \implies E\{\xi \mid \eta\} = E\xi.$

4. $\eta = f(\xi) \implies E\{\eta \mid \xi\} = E\{f(\xi) \mid \xi\} = f(\xi).$

5. $E(\eta f(\xi) \mid f(\xi)) = f(\xi)E\{\eta \mid \xi\}.$

6. $(\xi, \eta)^T \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \implies E(\xi \mid \eta) = a\eta + b.$

Замечание (Важное). Таким образом, если выборка нормальная, то зависимость линейная всегда.

7. $\operatorname{argmin}_{\hat{\eta} \in K = \{\phi(\xi)\}} E(\eta - \hat{\eta})^2 = E\{\eta \mid \xi\} = \hat{\eta}^*.$