

Статистика

Конспект практических занятий

Мат-Мех, ПМИ, СМ–СМ

\$Revision: ... (NEG: UNDER CONSTRUCTION) \$



Оглавление

I. Оценки характеристик и параметров распределения	6
1. Выборка и эмпирическая случайная величина	7
2. Виды признаков	8
3. Характеристики распределений и метод подстановки	9
4. Характеристики распределений и их оценки	11
4.1. Характеристики положения	11
4.2. Характеристики разброса	12
4.3. Характеристики формы распределения	13
4.4. Характеристики зависимости	14
5. Точечная оценка параметров распределения	15
5.1. Метод подстановки	15
5.2. Метод моментов	15
5.3. Метод оценки максимального правдоподобия	16
6. Свойства оценок	17
6.1. Несмещенность	17
6.2. Состоятельность	18
6.3. Асимптотическая нормальность	18
6.4. Эффективность	19
6.4.1. Эффективность и неравенство Рао-Крамера	19
6.5. Устойчивость оценок	20
II. Некоторые распределения, связанные с нормальным	21
1. Распределение $N(a, \sigma^2)$	22
2. Распределение $\chi^2(m)$	23
3. Распределение Стьюдента $t(m)$	24
4. Распределение Фишера	25
5. Квадратичные формы от нормально распределенных случайных величин	26
6. Распределение важных статистик	27
III. Проверка гипотез и доверительные интервалы	29
1. Построение критерия	31
1.1. Общие сведения	31
1.2. Схема построения критерия на основе статистики критерия	31

1.3. Ошибки первого и второго рода	32
1.4. Понятие вероятностного уровня p -value.	34
2. Проверка гипотезы о значении параметра (характеристики)	35
2.1. Проверка гипотезы о значении мат. ожидания (t -критерий)	35
2.1.1. $D\xi = \sigma^2 < \infty$	35
2.1.2. $D\xi$ неизвестна	35
2.1.3. Проверка гипотезы о мат.ож. в модели с одним параметром	35
2.2. Проверка гипотезы о значении дисперсии в нормальной модели (критерий χ^2)	36
2.2.1. $E\xi = a < \infty$	36
2.2.2. $E\xi$ неизвестно	36
2.3. Асимптотический критерий для гипотезы о значении параметра на основе MLE	36
3. Доверительные интервалы	38
3.1. Мотивация и определение	38
3.2. Доверительный интервал для проверки гипотезы о значении параметра	38
3.3. Доверительные интервалы для математического ожидания и дисперсии в нормальной модели	38
3.3.1. Доверительный интервал для a	38
3.3.2. Доверительный интервал для σ^2	39
3.4. Асимптотический доверительный интервал для математического ожидания в модели с конечной дисперсией	39
3.5. Асимптотический доверительный интервал для параметра на основе MLE	40
3.6. Использование SE для проверки гипотез и построения доверительных интервалов	41
4. Критерии проверки гипотезы о согласии с видом распределения	42
4.1. Критерий χ^2	42
4.1.1. Распределение с известными параметрами	42
4.1.2. Распределение с неизвестными параметрами	43
5. Критерий Колмогорова-Смирнова согласия с видом распределения	44
5.1. Произвольное абсолютно непрерывное распределение	44
6. Визуальное определение согласия с распределением	45
6.1. P-P plot	45
6.2. Q-Q plot	45
IV. Корреляционный анализ	46
1. Вероятностная независимость	48
1.1. Визуальное определение независимости	48
1.2. Критерий независимости χ^2	48
2. Линейная / нелинейная зависимость	50
2.1. Определение вида зависимости	50
2.2. Коэффициент корреляции Пирсона	50
2.2.1. Оценка коэффициента корреляции	51
2.2.2. Значимость коэффициента корреляции	51
2.3. Происхождение и сравнение других мер зависимости	51
2.3.1. Свойства корреляционного отношения	52
2.3.2. Выборочное корреляционное отношение	52
3. Частная корреляция	54

4. Зависимость между порядковыми признаками	56
4.1. Ранговый коэффициент Спирмана	56
4.1.1. Согласованность ρ и ρ_S	57
4.2. Ранговый коэффициент Кэндалла $\tau(\xi, \eta)$	58
V. Дисперсионный анализ	
NEG: не проверяла	60
1. Однофакторный дисперсионный анализ (One-way ANOVA ¹)	61
2. Множественные сравнения	63
2.1. Single	64
2.2. Stepdown (Holm's algorithm)	64
2.2.1. Частный случай	65
3. ANOVA Post-Hoc Comparison	66
3.1. Least Significant Difference (LSD)	66
3.2. Распределение размаха	66
3.3. Tukey's Honest Significant Difference (HSD) Test	67
3.4. Другие критерии	67
3.5. Scheffé's Method	68
3.6. Сравнение мощностей	68
VI. Регрессионный анализ	69
1. Регрессия	70
2. Парная линейная регрессия	71
2.1. Модель линейной регрессии	71
2.2. Доверительные интервалы для β_1 и β_2	72
3. Множественная линейная регрессия	74
3.1. Псевдо-обратные матрицы	74
3.2. Проекторы на подпространства	74
3.3. Ordinary and Total Least Squares	75
3.4. Свободный член	76
3.5. Стандартизованные признаки	76
3.6. Свойства оценки $\hat{\mathbf{b}}$	76
3.7. Свойства $\hat{\mathbf{b}}^{(c)}$ и $\hat{\mathbf{b}}^{(s)}$	77
3.8. Сравнение оценок	78
3.9. Разложение суммы квадратов и оценка σ^2	78
3.10. Проверка значимости коэффициентов линейной регрессии и доверительных интервалов	79
3.10.1. Расстояние Махаланобиса	79
3.10.2. Доверительный эллипсоид	79
3.11. Значимость регрессии	80
3.12. Анализ оценок коэффициентов	82
3.12.1. Корреляция между оценками коэффициентов в двумерном случае	82
3.12.2. Избыточность (redundancy) и ручное удаление признаков	82
3.12.3. Проверка гипотезы о том, что набор признаков избыточен	83

¹ANalysis Of VARIation

3.12.4. Stepwise автоматическое удаление/добавление признаков	83
3.12.5. Выбор модели на основе информационных критериев AIC и BIC	84
3.12.6. О множественном коэффициенте корреляции и саппрессорах	84
3.12.7. Как понять, что все хорошо	84
3.12.8. Заполнение пропусков	85
3.13. Анализ аутлаеров	85
3.13.1. Matrix plot	85
3.13.2. Deleted residuals	85
3.13.3. Studentized residuals	85
3.13.4. Расстояние по Куку и расстояние Махаланобиса	86
3.14. Проверка правильности и выбор модели	87
3.15. Доверительные интервалы для среднего предсказания и предсказательные интервалы	87
3.16. Сведение нелинейной модели к линейной	88
4. Модификации линейной регрессии.	89
4.1. Взвешенная регрессия (Weighted Least Squares)	89
4.2. Гребневая (Ridge) регрессия	90
A. Свойства условного математического ожидания	91

Часть I.

Оценки характеристик и параметров распределения

1. Выборка и эмпирическая случайная величина

Пусть $\xi : (\Omega, \mathcal{F}, \mathbb{P}) \rightarrow (V, \mathfrak{A})$ — случайная величина с распределением \mathcal{P} (пишем $\xi \sim \mathcal{P}$).

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} . (Знак \perp иногда используется для обозначения независимости случайных величин.)

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

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

Замечание. Подходящее определение выбирается по контексту.

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

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

Если ξ имеет дискретное распределение, то выборку можно *сгруппировать*; тогда распределение случайной величины $\hat{\xi}_n$ запишется как

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

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),*

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

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

Выше используется сходимость по вероятности (\mathbb{P} от слова probability). Обозначение $\zeta_n \xrightarrow{\mathbb{P}} \zeta$ означает, что $\forall \varepsilon > 0, \lim_{n \rightarrow \infty} \mathbb{P}(|\zeta_n - \zeta| > \varepsilon) = 0$.

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: $V \subset \mathbb{R}$, заданы операции с вещественными числами.

По типу операций:

- Аддитивные: заданы, т.е. имеют смысл в контексте данного признака, операции $+$, $-$. Разница между значениями характеризуется разностью значений.
- Мультипликативные: заданы операции \times , $/$; признак принимает не отрицательные значения. Разница между значениями измеряется в процентах (определяется делением).

По типу данных:

- Непрерывные
- Дискретные

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

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

Пример. Цвет глаз, имена, пол.

3. Характеристики распределений и метод подстановки

Определение. *Статистика* — измеримая функция от выборки.

Обобщением статистики является понятие характеристики.

Определение. *Характеристика* распределения — функционал от распределения:

$$T : \{\mathcal{P}\} \rightarrow D;$$

Чаще всего, $D = \mathbb{R}$.

Определение. *Оценка* — функция от выборки, не зависящая от генеральной характеристики θ .

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

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

Здесь используется обозначение cdf , сокращение от cumulative distribution function, т.е. от названия функции распределения. Мы здесь оставим это обозначение, однако часто на занятиях будем обозначать функцию распределения просто $F(x)$ или $F_\xi(x)$, чтобы подчеркнуть, какой случайной величины это функция распределения.

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

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

(сходимость a.s. означает сходимость almost surely, почти наверно, почти всегда; она вводится для случайных величин, заданных на вероятностном пространстве и означает, что сходимость имеет место для почти всех элементарных событий кроме, может быть, событий меры ноль).

Может возникнуть вопрос, почему слева от стрелки случайная величина. Дело в том, что в этом утверждении (как и в любом теоретическом утверждении в математической статистике, выборка понимается как «до эксперимента», т.е. все x_i — случайные величины, поэтому эмпирическая функция распределения в точке x , равная числу x_i , меньших x , тоже является случайной.

Более того, если cdf_ξ непрерывна, скорость сходимости имеет порядок $1/\sqrt{n}$ по теореме Колмогорова:

$$\sqrt{n} \sup_{x \in \mathbb{R}} \left| \widehat{\text{cdf}}_\xi(x) - \text{cdf}_\xi(x) \right| \xrightarrow{d} \mathcal{P}_{\text{K.S.}},$$

где $\mathcal{P}_{\text{K.S.}}$ — распределение Колмогорова-Смирнова.

Выше используется другой тип сходимости, более слабый, по распределению (d от слова distribution). Эта сходимость означает, что функция распределения случайной величины слева сходится в функции распределения, указанного справа, во всех точках ее непрерывности.

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

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

Следствие. *Значит, при достаточно больших n , в качестве интересующей характеристики $\theta = f(\xi)$ распределения \mathcal{P}_ξ можем брать ее оценку $\hat{\theta} = \hat{\theta}_n = f(\hat{\xi}_n)$ — аналогичную характеристику $\hat{\mathcal{P}}_n$. Этот метод называется методом подстановки.*

4. Характеристики распределений и их оценки

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

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

4.1. Характеристики положения

В качестве характеристики положения для количественных признаков выделяется 1-й момент — математическое ожидание и его оценка *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 (квантилью уровня p) называется

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

Квартиль есть квантиль уровня, кратного $1/4$; *дециль* — $1/10$; *перцентиль* — $1/100$. Эти характеристики определены для порядковых признаков (и, следовательно, для количественных тоже).

Замечание. \sup берется для учета случая не непрерывных функций распределения.

Определение. Медиана есть 0.5-квантиль:

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

Определение. Мода ($\text{mode } \xi$) есть точка локального максимума плотности или состояние с максимальной вероятностью для качественных признаков.

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:

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

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.

4.2. Характеристики разброса

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

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

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

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

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 m -элементного множества вычисляется по *формуле Хартли*:

$$H = \log_2 m$$

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

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

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

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

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.

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

SD позволяет оценивать стандартное отклонение распределения ξ .

Пусть $\hat{\theta}_n$ — статистика. Она имеет какое-то своё распределение, стандартное отклонение которого можно также оценить.

Определение. *Стандартная ошибка* оценки есть

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

Это показатель разброса оценки случайной величины.

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

$$\text{SE} := \text{SE}(\bar{x}) = \sqrt{\widehat{\text{D}\bar{x}}} = \sqrt{\frac{\widehat{\text{D}\xi}}{n}} = \frac{s}{\sqrt{n}}.$$

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

Пусть $c_\gamma = \text{qnt}_{N(0,1)} \gamma$. $N(\mu, \sigma^2)$ обозначает нормальное распределение с математическим ожиданием μ и дисперсией σ^2 .

Пример (С мостом и машинами). При возведении моста требуется, чтобы под ним могли проехать, условно, 95% машин. Чтобы эту высоту вычислить, достаточно собрать выборку высоты кузова проезжающих машин. Тогда нахождение искомой величины можно наглядно представить как выбор такой квантили гистограммы выборки, что суммирование соответствующих вероятностей даст $\gamma = 0.95$. В предположении, что выборка из нормального распределения, с более устойчивой оценкой квантили, интервал будет иметь вид

$$(\bar{x} \pm \text{SD} \cdot c_\gamma).$$

SE как показатель разброса выборочного среднего использовать по смыслу нельзя.

Пример (С паромом). Число машин, которое способен перевезти паром, есть Грузоподъемность/ $\mathbb{E}\xi$, где ξ — вес машины. Поскольку оценка \bar{x} всегда считается с погрешностью относительно истинного значения, интервал допустимого числа машин будет иметь вид

$$\frac{\text{Грузоподъемность}}{\bar{x} \pm \text{SE} \cdot c_\gamma}.$$

Подробности, включая определение c_γ см. в разделе посвященном доверительным интервалам (Глава 3).

4.3. Характеристики формы распределения

Для удобства, обозначим $\sigma^2 = m_2^{(0)} = \text{D}\xi$.

Определение. *Коэффициент асимметрии Пирсона* («скошенности»²)

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

Замечание. Не зависит от линейных преобразований.

Замечание. Старое определение скошенности было $\frac{\mathbb{E}\xi - \text{med}\xi}{\sigma}$.

Замечание. Типичный случай соответствует тому, что при положительном коэффициенте асимметрии «хвост вправо».

² «Skewness».

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 γ_4 для $N(0, 1)$.

Замечание. Положительный коэффициент эксцесса соответствует медленному убыванию на концах отрезка. Причём, так как распределение стандартизуется, имеется в виду убывание на хвостах, которое медленнее по порядку (!), чем убывание на хвостах у нормального распределения. Например, сравните e^{-x^2} , $e^{-x^2/10}$ и e^{-10x} . Часто говорят об островершинности при положительном эксцессе, но это просто вторая сторона скорости убывания на хвостах. Медленное убывание на хвостах означает на практике, что далекие от среднего значения встречаются необычно часто.

4.4. Характеристики зависимости

Определение. Пусть $(\xi_1, \xi_2)^T \sim \mathcal{P}$ и $(x_1, y_1)^T, \dots, (x_n, y_n)^T \sim \mathcal{P}$ — выборка из этого распределения. Тогда можно записать две другие важные характеристики: ковариацию и коэффициент корреляции:

$$\begin{aligned} \text{cov}(\xi_1, \xi_2) &= E(\xi_1 - E\xi_1)(\xi_2 - E\xi_2) = E\xi_1\xi_2 - E\xi_1E\xi_2 & \widehat{\text{cov}}(\mathbf{x}, \mathbf{y}) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ \text{cor}(\xi_1, \xi_2) &= \frac{\text{cov}(\xi_1, \xi_2)}{\sigma_{\xi_1}\sigma_{\xi_2}} & \widehat{\text{cor}}(\mathbf{x}, \mathbf{y}) &= \frac{\widehat{\text{cov}}(\mathbf{x}, \mathbf{y})}{s(\mathbf{x})s(\mathbf{y})}. \end{aligned}$$

5. Точечная оценка параметров распределения

5.1. Метод подстановки

Метод подстановки заключается в подстановке вместо неизвестного теоретического распределения известного эмпирического распределения. Например, вас интересует некоторая характеристика $f(\xi)$, а вы в качестве оценки предлагаете $\widehat{f(\xi)} = f(\xi_n)$, где $\hat{\xi}_n = \xi_n$ — эмпирическая случайная величина.

5.2. Метод моментов

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

$$\begin{cases} \mathbb{E}g_1(\xi) = \phi_1(\theta_1, \dots, \theta_r) \\ \vdots \\ \mathbb{E}g_r(\xi) = \phi_r(\theta_1, \dots, \theta_r) \end{cases} \implies \begin{cases} \theta_1 = f_1(\mathbb{E}g_1(\xi), \dots, \mathbb{E}g_r(\xi)) \\ \vdots \\ \theta_r = f_r(\mathbb{E}g_1(\xi), \dots, \mathbb{E}g_r(\xi)). \end{cases}$$

Примем

$$\theta_i^* = f_i(\hat{\mathbb{E}}g_1(\xi), \dots, \hat{\mathbb{E}}g_r(\xi)).$$

Часто, $g_i(\xi) = \xi^i$. Или, еще чаще, $g_1(\xi) = \xi$ и $g_i(\xi) = (\xi - \mathbb{E}\xi)^i$, $i > 1$, так как для таких моментов обычно известны формулы.

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 5.1 ($r = 1$). $\xi \sim U(0, \theta)$.

- Оценка по 1-му моменту: $g(\xi) = \xi$ и

$$\mathbb{E}\xi = \int_0^\theta \frac{1}{\theta} x \, dx = \frac{1}{\theta} \frac{x^2}{2} \Big|_0^\theta = \frac{\theta}{2} \implies \theta = 2\mathbb{E}\xi, \quad \theta^* = 2\bar{x}.$$

- Оценка по k -му моменту: $g(\xi) = \xi^k$ и

$$\mathbb{E}\xi^k = \frac{1}{\theta} \int_0^\theta x^k \, dx = \frac{1}{\theta} \frac{x^{k+1}}{k+1} \Big|_0^\theta = \frac{\theta^k}{k+1} \implies \theta^* = \sqrt[k]{(k+1) \frac{1}{n} \sum_{i=1}^n x_i^k}.$$

Пример 5.2 ($r = 1$). Пусть $\xi \sim \text{Exp}(\lambda)$. Тогда $\mathbb{E}\xi = \lambda$ и $\bar{x} = \lambda$.

5.3. Метод оценки максимального правдоподобия

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

Определение. Пусть

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

Тогда *функция правдоподобия* определяется как значение распределения выборки (плотности в непрерывном случае и вероятности значений в дискретном) с подстановкой выборки вместо аргумента:

$$L(\boldsymbol{\theta} | \mathbf{x}) = P(\mathbf{x} | \boldsymbol{\theta}).$$

Пример 5.3. Пусть $\xi \sim N(\mu, \sigma^2)$. По независимости x_i , $p_{\boldsymbol{\theta}}(\mathbf{x})$ распадается в произведение:

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

Пример 5.4. $\xi \sim \text{Pois}(\lambda)$,

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

Утверждение. Пусть \mathbf{x} — выборка. В качестве оценки максимального правдоподобия¹ $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ следует взять

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

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

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

откуда

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

Утверждение. В условиях регулярности:

1. Существует один глобальный максимум, так что

$$\left. \frac{\partial \ln L(\lambda | \mathbf{x})}{\partial \lambda} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{MLE}}} = 0.$$

2. $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ обладает всеми свойствами (про определение этих свойств написано в следующих разделах):

- а) Состоятельность;
- б) Асимптотическая несмещенность;
- в) Асимптотическая нормальность;
- г) Асимптотическая эффективность.

¹Maximum likelihood estimate (MLE).

6. Свойства оценок

6.1. Несмещенность

Определение. *Смещение*¹ есть

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

Определение. *Среднеквадратичная ошибка*² есть

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

Замечание. Поскольку

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

то

$$\underbrace{E(\hat{\theta}_n - \theta)^2}_{\text{MSE}} = D\hat{\theta}_n + \underbrace{(E(\hat{\theta}_n - \theta))^2}_{\text{bias}^2}. \quad (6.1)$$

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

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

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

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

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

□

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

Доказательство. Для $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ имеем:

$$\begin{aligned} Es^2 &= \frac{1}{n} \sum_{i=1}^n E(x_i - E\xi)^2 - E(\bar{x} - E\xi)^2 = \frac{1}{n} \sum_{i=1}^n Dx_i - D\bar{x} = D\xi - \frac{1}{n} D\xi \\ &= \frac{n-1}{n} D\xi \xrightarrow{n \rightarrow \infty} D\xi. \end{aligned}$$

□

Определение. *Исправленная дисперсия:*

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

Очевидно, исправленная дисперсия — несмещенная оценка дисперсии.

¹Bias.

²Mean squared error (MSE).

6.2. Состоятельность

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

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

Как следует из равенства (6.1), для асимптотически несмещенных оценок состоятельность в средне-квадратическом следует из сходимости дисперсии оценки к нулю.

Определение. Оценка называется *состоятельной*, если

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

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

Доказательство. В самом деле, по неравенству Чебышева,

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

□

Предложение. \hat{m}_k является состоятельной оценкой m_k .

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

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \xrightarrow{P} m_1(\mathcal{P}).$$

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

□

Замечание. Состоятельность выполняется и для центральных моментов $m_k^{(0)}$, так как они выражаются через нецентральные, а свойство состоятельности сохраняется для линейной комбинации.

В частности, \bar{x} — состоятельная оценка $E\xi$ и s^2 — состоятельная оценка $D\xi$.

6.3. Асимптотическая нормальность

Определение. Оценка $\hat{\theta}_n$ называется *асимптотически нормальной* оценкой параметра θ с коэффициентом $\sigma^2(\theta)$ если

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \sigma^2(\theta)).$$

Пример. \bar{x} — асимптотически нормальная оценка, если $D\xi < \infty$, $D\xi \neq 0$:

$$\sqrt{n}(\bar{x} - E\xi) \xrightarrow{d} N(0, D\xi).$$

Доказательство. По ЦПТ,

$$\sqrt{n}(\bar{x} - E\xi) = \frac{\sum_{i=1}^n x_i - nE\xi}{\sqrt{n}} \xrightarrow{d} N(0, D\xi).$$

□

Мы обсуждали, что асимптотическую нормальность можно определять и в слабом смысле — как сходимость по распределению к нормальному распределению $N(0, 1)$ стандартизированной случайной величины.

6.4. Эффективность

Определение. Говорят, что оценка $\hat{\theta}^{(1)}$ лучше $\hat{\theta}^{(2)}$ в среднеквадратичном смысле, если

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

Замечание. Для несмещенных оценок определение эквивалентно, конечно,

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

Определение. В классе несмещенных оценок оценка называется эффективной (в средне-квадратическом), если ее дисперсия минимальна. В классе асимптотически несмещенных оценок оценка $\hat{\theta} = \hat{\theta}_n$ называется асимптотически эффективной, если для любой другой оценки $\hat{\theta}^*$ выполнено $\lim_{n \rightarrow \infty} \text{D}\hat{\theta}_n / \text{D}\hat{\theta}_n^* \leq 1$.

6.4.1. Эффективность и неравенство Рао-Крамера

Пусть $\mathcal{P}_\xi(\theta)$, $\theta = (\theta_1, \dots, \theta_r)^\top$ — параметрическая модель. Пусть $r = 1$.

Определение. Информанта n -го порядка:

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

Определение. Информационное количество Фишера:

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

Утверждение.

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

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

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

Замечание.

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

откуда, для повторной независимой выборки,

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

Определение. $C \subset \mathbb{R}$ есть носитель параметрического семейства распределений $\mathcal{P}(\theta)$, если

$$\xi \sim \mathcal{P}(\theta) \implies \text{P}(\xi \in C) = 1, \quad \forall \theta \in \Theta.$$

Определение. Условие регулярности: имеют отношение к независимости носителя распределения от параметра, а также к существованию и ограниченности производных функции лог-правдоподобия по параметру до определённого порядка дифференцирования.

Пример. $\text{Exp}(\lambda)$ — регулярное семейство; $\text{U}(0, \theta)$ — не является регулярным.

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

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

Для смещенных оценок,

$$\text{D}\hat{\theta}_n \geq \frac{(1 + \text{bias}'(\theta))^2}{I_n(\theta)}.$$

Следствие. Несмещенная оценка является эффективной, если:

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

Следствие. Асимптотически несмещенная оценка является асимптотически эффективной, если:

$$D\hat{\theta}_n \cdot I_n(\theta) \rightarrow 1 \text{ as } n \rightarrow \infty.$$

Упражнение (Хорошее). Показать, что \bar{x} является эффективной оценкой μ в модели $\xi \sim N(\mu, \sigma^2)$.

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

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

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

то $\hat{\lambda}_n$ — эффективная оценка (по свойствам $\hat{\theta}_{\text{MLE}}$, гарантировано, что она асимптотически эффективная).

6.5. Устойчивость оценок

Так как в реальных данных часто бывают те или иные ошибки, часто жертвуют точностью для увеличения устойчивости (робастности) к выбросам. Устойчивые аналоги оценок часто строятся на основе рангов (номеров по порядку в упорядоченной выборке). Приведем пример.

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.

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

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

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

откуда

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

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

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

Часть II.

Некоторые распределения, связанные с нормальным

1. Распределение $N(a, \sigma^2)$

Свойства хорошо известны. В частности, плотность имеет вид

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

математическое ожидание равно a , дисперсия σ^2 , асимметрия и эксцесс равны 0.

Рассмотрим вопрос про измерение расстояния в сигмах. Будет говорить, что точка далеко от мат.ожидания, если это и более далекие значения маловероятны.

Формально, пусть $\xi \sim N(a, \sigma^2)$. Рассмотрим $P(|\xi - a| > k\sigma)$. Эта вероятность не зависит от σ и равна $2(1 - \Phi(k))$, где $\Phi(x)$ — функция стандартного нормального распределения $N(0, 1)$.

Значения $P(|\xi - a| > k\sigma)$:

k	Вероятность
1	0.317
1.64	0.101
1.96	0.050
2	0.046
3	2.70E-03
6	1.97E-09

Отсюда правило двух сигма (вероятность быть на расстоянии от мат.ож. больше двух сигм примерно равна 0.05), правило трех сигм, правило шести сигм.

2. Распределение $\chi^2(m)$

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

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

Свойства¹ $\chi^2(m)$

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

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

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

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

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

и

$$\text{cdf}_{\chi^2(50)}(80) = 0.9955 \approx \Phi(3) = 0.9986.$$

Предложение. $\chi^2(m)/m \xrightarrow{m \rightarrow \infty} 1$.

Доказательство. По ЗБЧ. □

¹Вычисление $D\eta$: <https://www.statlect.com/probability-distributions/chi-square-distribution>

3. Распределение Стьюдента $t(m)$

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

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

Свойства $t(m)$

- При $m = 1$ это распределение Коши, у него не существует математического ожидания.
- При $m > 1$, $E\xi = 0$ по симметричности.
- При $m > 2$, $D\xi = m/(m - 2)$.
- При $m > 3$, $A\xi = 0$ по симметричности.
- При $m > 4$, $K\xi = 6/(m - 4)$.

Предложение. *Распределение Стьюдента сходится к стандартному нормальному:*

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

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

□

4. Распределение Фишера

Определение. Распределение Фишера имеет вид

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

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

5. Квадратичные формы от нормально распределенных случайных величин

(Это на след. семестр, сейчас можно не вникать.)

Пусть $\boldsymbol{\xi} = (\xi_1, \dots, \xi_p)^\top \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_p)$, \mathbf{B} — симметричная, неотрицательно определенная матрица. Найдем распределение $\boldsymbol{\xi}^\top \mathbf{B} \boldsymbol{\xi}$.

Утверждение. Пусть $\boldsymbol{\xi} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_p)$, \mathbf{B}, \mathbf{C} — симметричные матрицы размерности $p \times p$. Тогда $\boldsymbol{\xi}^\top \mathbf{B} \boldsymbol{\xi} \perp \boldsymbol{\xi}^\top \mathbf{C} \boldsymbol{\xi} \iff \mathbf{BC} = \mathbf{0}$.

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

$$\begin{aligned}\bar{x}^2 &= \frac{1}{n^2} \left(\sum_{i=1}^n x_i \right)^2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n x_i x_j = \frac{1}{n} \mathbf{x} \underbrace{\begin{pmatrix} 1/n & \dots & 1/n \\ \vdots & \ddots & \vdots \\ 1/n & \dots & 1/n \end{pmatrix}}_{\mathbf{B}} \mathbf{x}^\top \\ s^2 &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n} \mathbf{x} \mathbf{B} \mathbf{x}^\top = \frac{1}{n} \left(\mathbf{x} \mathbf{I}_n \mathbf{x}^\top - \mathbf{x} \mathbf{B} \mathbf{x}^\top \right) = \frac{1}{n} \mathbf{x} \underbrace{\begin{pmatrix} 1 - 1/n & \dots & -1/n \\ \vdots & \ddots & \vdots \\ -1/n & \dots & 1 - 1/n \end{pmatrix}}_{\mathbf{C} = \mathbf{I}_n - \mathbf{B}} \mathbf{x}^\top.\end{aligned}$$

Таким образом, $n\bar{x}^2 = \mathbf{x} \mathbf{B} \mathbf{x}^\top$ и $ns^2 = \mathbf{x} \mathbf{C} \mathbf{x}^\top$. Но

$$\mathbf{BC} = \mathbf{B}(\mathbf{I}_n - \mathbf{B}) = \mathbf{B} - \mathbf{B}^2 = \mathbf{0},$$

так как

$$\mathbf{B}^2 = \begin{pmatrix} 1/n & \dots & 1/n \\ \vdots & \ddots & \vdots \\ 1/n & \dots & 1/n \end{pmatrix}^2 = \begin{pmatrix} n \cdot 1/n & \dots & n \cdot 1/n \\ \vdots & \ddots & \vdots \\ n \cdot 1/n & \dots & n \cdot 1/n \end{pmatrix} = \mathbf{B}.$$

Значит, $\bar{x}^2 \perp s^2$.

Видно, что $\sigma^{-2} \boldsymbol{\xi}^\top \mathbf{I}_p \boldsymbol{\xi} \sim \chi^2(p)$. 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

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

Пример. Покажем, что

$$ns^2/\sigma^2 \sim \chi^2(p-1).$$

Воспользуемся представлением из предыдущего примера: $ps^2 = \mathbf{x}^\top \mathbf{C} \mathbf{x}$. 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 $\boldsymbol{\xi} \sim N(\mathbf{0}, \mathbf{I}_p)$, $\boldsymbol{\xi}^\top \boldsymbol{\xi} = \sum_i Q_i$, 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. $Q_i \perp Q_j, \quad \forall i \neq j$, т.е. $\mathbf{B}_i \mathbf{B}_j = \mathbf{0}$.

6. Распределение важных статистик

Пусть $\xi \sim N(a, \sigma^2)$.

Предложение. $t = \sqrt{n} \frac{(\bar{x} - E\xi)}{\sigma}$ имеет стандартное нормальное распределение.

Доказательство.

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

□

Определим $s_a^2 = \sum_{i=1}^n (x_i - a)^2 / n$.

Предложение. $ns_a^2 / \sigma^2 \sim \chi^2(n)$.

Доказательство.

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

□

Предложение. $ns^2 / \sigma^2 = (n-1)\hat{s}^2 / \sigma^2 \sim \chi^2(n-1)$.

Доказательство. См. раздел 5).

□

Альтернативное доказательство. По определению запишем

$$\underbrace{D\hat{\xi}_n}_{s^2} = D(\hat{\xi}_n - a) = \underbrace{E(\hat{\xi}_n - a)^2}_{s_a^2} - \underbrace{(E(\hat{\xi}_n - a))^2}_{(\bar{x} - a)^2}.$$

Домножив обе части на n/σ^2 , получим

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

□

Замечание. Для строгого доказательства, нужно использовать независимость \bar{x}^2 и s^2 (см. раздел 5).

Предложение. Следующая статистика имеет распределение Стьюдента:

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

Предложение. $t = \sqrt{n-1} \frac{\bar{x}-a}{s} = \sqrt{n} \frac{\bar{x}-a}{\tilde{s}} \sim t(n-1)$.

Доказательство.

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

поскольку

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

и они независимы (также пока без доказательства — используется техника квадратичных форм или можно доказать через разложение дисперсии). \square

Часть III.

Проверка гипотез и доверительные интервалы

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

1. Построение критерия

1.1. Общие сведения

Пусть H_0 — это гипотеза, т.е. некоторое предположение о случайной величине ξ , которое мы хотим проверить (модель — это предположение, которое считается верным без проверки).

Задаем уровень значимости $0 < \alpha < 1$.

Тогда критерий — это разбиение множества V^n всевозможных значений выборки \mathbf{x} на две области, критическую $A_\alpha^{(\text{крит})}$ и доверительную $A_\alpha^{(\text{дов})}$ так, что $\alpha_I = P_{H_0}(\mathbf{x} \in A_\alpha^{(\text{крит})}) = \alpha$ (P_{H_0} — вероятность, соответствующая предположению, что H_0 верна).

При проверки гипотезы (уже после эксперимента), если выборка попала в критическую область $A_\alpha^{(\text{крит})}$, то нулевая гипотеза отвергается, а если в доверительную $A_\alpha^{(\text{дов})}$, то не отвергается (важно, что именно нет основания отвергнуть, но и принять нельзя).

Обычно разбиение строят с помощью статистики критерия $t = t(\mathbf{x})$. В этом случае на доверительную и критическую область нужно делать область значений статистики критерия, а это подмножество вещественных чисел. Поэтому критическая область $A_\alpha^{(\text{крит})}$ выбирается так, чтобы $\alpha_I = P_{H_0}(t \in A_\alpha^{(\text{крит})}) = \alpha$.

Допустимо строить разбиение так, чтобы выполнялось $\alpha_I \leq \alpha$ (тогда критерий называется консервативным).

Часто удается построить только асимптотический критерий, когда $\alpha_I \rightarrow \alpha$ при $n \rightarrow \infty$. В этом случае критерий можно применять при достаточно (для критерия) большом объеме выборки, где допустимый объем выборки зависит от скорости сходимости.

Ниже более подробно.

1.2. Схема построение критерия на основе статистики критерия

1. Строим статистику критерия t так, что:

- Статистика критерия t должна измерять то, насколько выборка соответствует гипотезе. В этом случае мы получаем значение статистики критерия для «идеального соответствия».

Например, если гипотеза про математическое ожидание, то $t = \bar{x} - E\xi$ подходит под это требование. Если гипотеза про дисперсию, то соответствие правильнее измерять отношением и поэтому подошло бы $t = s^2/D\xi$.

Пример. Пусть $H_0 : E\xi = a_0$; тогда $t = \bar{x} - a_0$ и «идеальное значение» $t = 0$.

- Распределение t при верной H_0 должно быть известно хотя бы асимптотически. Из-за этого часто преобразовывают вариант меры несоответствия, приведенный выше.

Пример. См. раздел 2.1.

2. Строим разбиение области значений статистики критерия t так, что:

- $P(t \in A_\alpha^{(\text{крит})}) = \alpha$.
- Если альтернативная гипотеза H_1 не конкретизирована, то $A_\alpha^{(\text{крит})}$ следует выбрать так, чтобы она располагалась как можно дальше от идеального значения.

Пример. В случае $t \sim N(0, 1)$ при идеальном значении 0, разумно определить $A_\alpha^{(\text{крит})}$ «на хвостах» графика $\text{pdf}_{N(0,1)}$ симметрично по обе стороны от 0 так, что для $A_\alpha^{(\text{крит})} =$

1. Построение критерия

$$(-\infty, -t_\alpha) \cup (t_\alpha, \infty)$$

$$\alpha/2 = \int_{-\infty}^{-t_\alpha} \text{pdf}_{N(0,1)}(y) dy = \int_{t_\alpha}^{+\infty} \text{pdf}_{N(0,1)}(y) dy.$$

Иными словами,

$$\alpha/2 = 1 - \text{cdf}_{N(0,1)}(t_1) \implies t_1 = \text{cdf}_{N(0,1)}^{-1}(1 - \alpha/2)$$

и аналогично для t_0 .

- Если H_1 известна, то $A_\alpha^{(\text{крит})}$ выбирается так, чтобы максимизировать мощность критерия против альтернативы H_1 , определения см. ниже.

1.3. Ошибки первого и второго рода

Определение (Ошибки I-го и II-го родов). Пусть мы проверяем нулевую гипотезу H_0 . Зафиксируем альтернативную гипотезу H_1 — такое отклонение от H_0 , что его обнаружение важно для нас. Тогда

- ошибка I-го рода есть отвержение H_0 , при верной H_0 ; соответствующая вероятность есть

$$\alpha_I := P_{H_0}(\mathbf{x} \in A_{\text{крит}}^{(\alpha)}).$$

Замечание. Для точного критерия вероятность α_I совпадает с уровнем значимости α .

- ошибка II-го рода есть не отвержение H_0 при верной H_1 ; соответствующая вероятность есть

$$\alpha_{II} := P_{H_1}(\mathbf{x} \in A_{\text{дов}}^{(\alpha)}).$$

Определение. *Мощность* критерия против альтернативы есть

$$\beta := 1 - \alpha_{II} = 1 - P_{H_1}(\mathbf{x} \in A_{\text{дов}}^{(\alpha)}) = P_{H_1}(\mathbf{x} \in A_{\text{крит}}^{(\alpha)}).$$

Иными словами, это способность критерия отличать H_1 от H_0 .

Определение. Критерий называется *состоятельным*, если $\beta \rightarrow 1$.

Замечание. Утверждать об *отвержении* гипотезы можно с вероятностью ошибки α (достаточно малой); утверждать о *принятии* гипотезы можно с вероятностью ошибки α_{II} — не контролируемой и могущей быть довольно большой. Поэтому гипотезу H_0 можно только отвергать или не отвергать, так как мы контролируем ошибку неправильного решения (отвергнуть). Принимать гипотезу нельзя, так как мы не контролируем, вообще говоря, ошибку неправильного решения в случае принятия гипотезы.

Замечание. В связи с введенным понятием мощности можно описать две проблемы:

- проблема маленьких объемов выборки состоит в том, что в этом случае мощность маленькая и критерий не заметит отличие от H_1 от H_0 , т.е. с большой вероятностью не отвергнет H_0 , хотя будет верна H_1 ;
- как ни странно, но есть также проблема слишком больших объемов выборки, когда мощность слишком большая, т.е. критерий может отвергнуть H_0 с большой вероятностью, даже если она «почти» верна (например, из-за ошибок округления).

Заметим, что вероятность ошибки первого рода α_I фиксирована (как минимум, она ограничена сверху выбранным значением α , в то время как $\alpha_{II} = \alpha_{II}(\alpha, n, H_1)$). Ниже продемонстрируем эту зависимость на примере.

1. Построение критерия

Пример. Пусть $\xi \sim N(a, \sigma^2)$, σ^2 известна — это модель. Гипотезы имеют вид $H_0 : a = a_0$, $H_1 : a = a_1$. Тогда

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

В то же время, поскольку при верной H_1 , $E\bar{\mathbf{x}} = 1/n \cdot \sum_{i=1}^n \xi_i = n/n \cdot a_1$ то

$$Et = \frac{\sqrt{n}(a_1 - a_0)}{\sigma} \Rightarrow t \sim N\left(\frac{\sqrt{n}(a_1 - a_0)}{\sigma}, 1\right) \text{ при верной } H_1.$$

(дисперсия, конечно, не меняется при сдвиге).

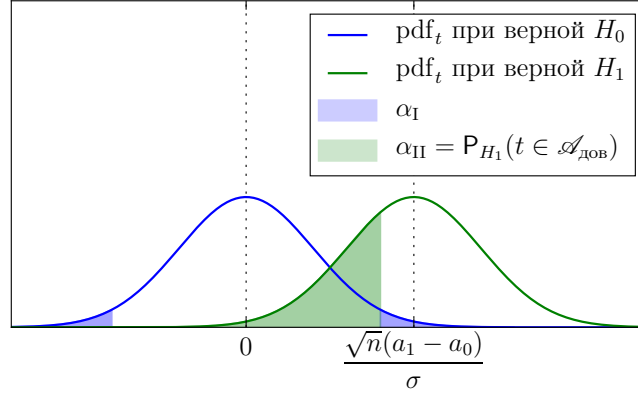


Рис. 1.1.: Плотности распределения t (неоптимальное разбиение)

Чтобы минимизировать α_{II} , логично определить $A_{\text{крит}}$ только на одном хвосте — с той стороны, где находится альтернатива.

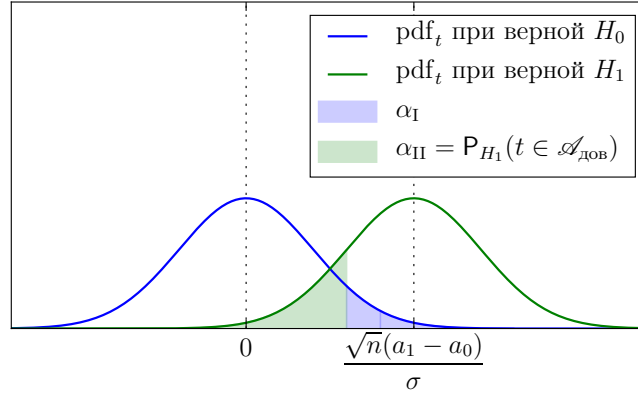


Рис. 1.2.: Плотности распределения t (оптимальное разбиение)

Таким образом, мы увидели, что если известна альтернатива, то можно выбирать критерий (= разбиение на доверительную и критическую области; статистика критерия — лишь удобной средство для этого), более мощный против именно этой альтернативы. И вообще, разные критерии для одной и той же основной гипотезы сравниваются по мощности против интересующей нас альтернативы.

Помимо этого, по рисунку видно, что ошибка второго рода α_{II} уменьшается при (1) увеличении ошибки первого рода (уровня значимости α) (поэтому мы ее выбираем максимальной из всех допустимых); (2) увеличении объема выборки (если α_{II} стремится к нулю при увеличении объема выборки, то критерий называется состоятельным после данной альтернативы) и (3) увеличении разницы между H_0 и H_1 в том же смысле, в каком статистика критерия измеряет разницу между

выборкой и нулевой гипотезой. Напомним, что мощность критерия против некоторой альтернативы (чего-то отличного от основной гипотезы) говорит нам о том, насколько хорошо критерий обнаруживает это отличие.

1.4. Понятие вероятностного уровня p -value.

Определение. p -value есть такое значение, что при значениях уровня значимости α , больших p -value, H_0 отвергается (по причине попадания t в $A_\alpha^{\text{крит}}$), а при меньших — не отвергается.

p -value — не вероятность, это пороговое значение. Неформально его можно интерпретировать как меру согласованности H_0 и выборки. Например, при больших значениях p -value практически при всех разумных уровнях значимости гипотеза не отвергается. При близких к нулю значениях p -value, наоборот, гипотеза будет отвергаться.

p -value — максимальное значение уровня значимости, при котором гипотеза не отвергается (значение статистики критерия попадает в доверит. область). Или, что эквивалентно, минимальное значение уровня значимости, при котором гипотеза отвергается.

Заметим, что для вычисления значения p -value даже для одной и той же статистики критерия нет единой формулы.

Неформальный комментарий, как находить p -value. Нужно сначала 1) нарисовать плотность статистики критерия, если верна H_0 , и на ней сделать разметку — где там критическая область, где доверительная для некоторого уровня значимости α (α равняется вероятности попасть в критическую область, т.е. равно площади под соотв. частью графика плотности). Разбиение на доверительную и критическую область обычно делают так, чтобы в критическую область попали значения, наиболее далекие от того значения статистики критерия, которое соответствует идеальному соотношению выборки и гипотезы. 2) нарисовать еще раз эту плотность и нанести значение статистики критерия.

Глядя на обе картинки, мысленно или рисуя, меняйте α и двигайте границу между доверительной и критической областями. Как упоминали выше, p -value — максимальное значение уровня значимости, при котором гипотеза не отвергается (значение статистики критерия попадает в доверит. область). Или, что эквивалентно, минимальное значение уровня значимости, при котором гипотеза отвергается. Замечу, что на основе этого определения сразу понятно, как, получив значение p -value, определить, при каких уровнях значимости (можно написать неравенство) гипотеза отвергается, а при каких — нет. Именно в таком виде и надо выдавать ответ.

На основе этих манипуляций должно стать понятно, что нужно для вычисления p -value — считать функцию распределения $F(t)$ или считать $1 - F(t)$, умножать потом на два или нет. После того как поймете, можно уже находить p -value с помощью R или Python. Например, в R есть `pnorm` — функция распределения нормального распределения, `pt` — распределения Стьюдента, `pf` — распределение Фишера.

На всякий случай, замечание про уровень значимости α : уровень значимости задается заранее, еще до проверки гипотезы. Его смысл — максимальная вероятность ошибки (отвергнуть H_0 неправильно), на которую согласен тот, кто будет отвечать за последствия ошибки. Но проблема в том, что те, кто делают обработку данных (в том числе, компьютеры) — это не те, которые отвечают за последствия, т.е. не те, кто устанавливают уровень значимости. Поэтому надо дать ответ в общем виде — в таком диапазоне уровней значимости гипотеза отвергается, а в таком — нет. Чтобы это сделать, все статистические пакеты выдают p -value. На его основе уже можно сформулировать общий ответ: при таких-то уровнях значимости гипотеза отвергается, при таких-то не отвергается. Иногда, когда некого спросить об уровне значимости, используются некоторые стандартные значения типа 0.05.

2. Проверка гипотезы о значении параметра (характеристики)

2.1. 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.

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

Предложение. Пусть $D\xi = \sigma^2 < \infty$; тогда используется следующая статистика (z -score)

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

Предложение. При условии $\xi \sim N(a, \sigma^2)$,

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

Доказательство.

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

□

2.1.2. $D\xi$ неизвестна

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

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

Сходимость к нормальному распределению следует из модифицированной теоремы Леви (модифицированная ЦПТ), которая позволяет заменять дисперсию на ее состоятельную оценку с сохранением сходимости к тому же нормальному распределению.

Предложение. При условии нормальности данных,

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

2.1.3. Проверка гипотезы о мат.ож. в модели с одним параметром

Разница с общим случаем состоит в том, что в параметрической модели с одним параметром не нужно оценивать дисперсию. Так как все выражается через этот параметр, то имеем формулу для дисперсии через значение параметра, предполагаемое в нулевой гипотезе.

z -критерий для пропорции в модели Бернулли Пусть $\xi \sim \text{Ber}(p)$. Поскольку $E\xi = p$, можно воспользоваться только что введенной статистикой; учитывая $D\xi = p(1-p)$, получаем статистику критерия для гипотезы $H_0 : p = p_0$:

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

2. Проверка гипотезы о значении параметра (характеристики)

z-критерий для интенсивности потока в модели Пуассона Пусть $\xi \sim \text{Pois}(\lambda)$. Поскольку $E\xi = \lambda$, можно воспользоваться только что введенной статистикой; учитывая $D\xi = \lambda$, получаем статистику критерия для гипотезы $H_0 : \lambda = \lambda_0$:

$$t = \sqrt{n} \frac{\bar{x} - \lambda_0}{\sqrt{\lambda_0}} \xrightarrow{d} N(0, 1).$$

2.2. 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.

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

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

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

2.2.2. $E\xi$ неизвестно

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

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

Упражнение. $s^2 = 1.44, \bar{x} = 55, n = 101$. 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:

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

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

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

┘

2.3. Асимптотический критерий для гипотезы о значении параметра на основе MLE

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

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

2. Проверка гипотезы о значении параметра (характеристики)

по асимптотической несмещенности,

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

и, учитывая асимптотическую эффективность ($D\hat{\theta}_{\text{MLE}} I_n(\theta) \xrightarrow[n \rightarrow \infty]{} 1$), запишем статистику для $H_0 : \theta = \theta_0$:

$$t = \left(\hat{\theta}_{\text{MLE}} - \theta_0 \right) \sqrt{I_n(\theta_0)} \xrightarrow{d} N(0, 1).$$

Задание: построить критерий для гипотезы о значении параметра для распределений Бернулли и Пуассона.

3. Доверительные интервалы

3.1. Мотивация и определение

Точечные оценки не дают информации о том, насколько (количественно) настоящее значение далеко от оценки.

Определение. $[b_1, b_2]$ — *доверительный интервал* для параметра θ с уровнем доверия $\gamma \in [0, 1]$, если $\forall \theta$

$$P(\theta \in [b_1, b_2]) = \gamma, \quad \text{где } b_1 = b_1(\mathbf{x}), b_2 = b_2(\mathbf{x}),$$

т.е. границы доверительного интервала — это статистики (функции от выборки, случайные величины «до эксперимента»).

Замечание. Если выборка из дискретного распределения, то b_1, b_2 — тоже дискретны. Поэтому наперед заданную точность получить может не получиться; в таких случаях знак « $=$ » заменяют « \geq ». Аналогично с заменой на « $\xrightarrow{n \rightarrow \infty}$ » для асимптотических доверительных интервалов, когда точные получить невозможно или затруднительно.

3.2. 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

$$P(\theta \in [b_1(\mathbf{x}), b_2(\mathbf{x})]) = \gamma = 1 - \alpha \implies \alpha = 1 - P(\theta \in [a_\gamma(\mathbf{x}), b_\gamma(\mathbf{x})]) = P(\theta \notin [a_\gamma(\mathbf{x}), b_\gamma(\mathbf{x})]).$$

Соответственно,

$$\begin{cases} \text{отвергаем } H_0, & \text{если } \theta_0 \notin [b_1(\mathbf{x}), b_2(\mathbf{x})] \\ \text{не отвергаем } H_0, & \text{если } \theta_0 \in [b_1(\mathbf{x}), b_2(\mathbf{x})]. \end{cases}$$

Вероятность ошибки первого рода равна α , что соответствует определению критерия. Заметим, что здесь мы пользуемся общим определением критерия, а не частным случаем, когда критерий строится через статистику критерия.

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

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

3.3.1. Доверительный интервал для a

- Пусть σ^2 известно. Свяжем a с выборкой через статистику критерия $t = \sqrt{n} \frac{(\bar{x} - a)}{\sigma} \sim N(0, 1)$:

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

3. Доверительные интервалы

Решений уравнения $P(c_1 < \sqrt{n}(\bar{x} - a)/\sigma < c_2) = \Phi(c_2) - \Phi(c_1) = \gamma$ 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,

$$P \left(a \in \left(\bar{x} \pm \frac{\sigma}{\sqrt{n}} c_\gamma \right) \right) = \gamma.$$

- Пусть σ^2 неизвестно. По аналогии,

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

и

$$P \left(a \in \left(\bar{x} \pm \frac{\tilde{s}}{\sqrt{n}} c_\gamma \right) \right) = \gamma.$$

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

Решение.

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

┘

3.3.2. 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

$$P \left(c_1 < \frac{ns_a^2}{\sigma^2} < c_2 \right) = \gamma \iff P \left(\sigma^2 \in \left(\frac{ns_a^2}{c_2}, \frac{ns_a^2}{c_1} \right) \right) = \gamma.$$

- Пусть a неизвестно. Тогда аналогично

$$P \left(\sigma^2 \in \left(\frac{ns^2}{c_2}, \frac{ns^2}{c_1} \right) \right) = \gamma,$$

где

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

3.4. 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 $E\xi = a$, не задавая параметрическую модель. Пусть $\{x_i\}$ i.i.d., тогда

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

Если заменить σ на ее состоятельную оценку (s), то по модифицированной теореме Леви (будет у Владимира Викторовича в след. году) сходимость не испортится. Тогда

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

Все это так же, как было в разделе 2.1.2.

Аналогично разделу 2.1.3, доверительные интервалы можно улучшить (сделать большее точными, т.е. вероятность попадания в них ближе к γ при фиксированном n) в параметрической модели, если вместо независимой оценки дисперсии использовать оценку, полученную на основе оценок параметров. Разница в том, что там можно было использовать значение параметра, взятое из гипотезы, а в доверительных интервалах придется подставлять оценки (или решать нелинейные неравенства, см. следующий раздел).

3.5. Асимптотический доверительный интервал для параметра на основе MLE

В точности, как было при проверке гипотез,

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

Чтобы по аналогии с предыдущим выразить θ в $P(c_1 < T < c_2) = P(|T| < c_\gamma) = \gamma$, необходимо знать зависимость $I_n(\theta)$ от θ . Для Pois и Ber разрешение неравенства относительно θ эквивалентно решению неравенства для квадратичного полинома.

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

$$P\left(-c_\gamma < (\hat{\theta}_{\text{MLE}} - \theta) \sqrt{I_n(\hat{\theta}_{\text{MLE}})} < c_\gamma\right) \rightarrow \gamma \iff \quad (3.1)$$

$$P\left(\theta \in \left(\hat{\theta}_{\text{MLE}} \pm \frac{c_\gamma}{\sqrt{I_n(\hat{\theta}_{\text{MLE}})}}\right)\right) \rightarrow \gamma, \quad (3.2)$$

где

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

Пример. $\xi \sim \text{Pois}(\lambda)$. В разделе 5.3 получали: $\hat{\lambda}_{\text{MLE}} = \bar{x}$ и $I_n(\lambda) = n/\lambda$ и $I_n(\hat{\lambda}) = n/\bar{x}$, откуда

$$P\left(\lambda \in \left(\bar{x} \pm c_\gamma \frac{\sqrt{\bar{x}}}{\sqrt{n}}\right)\right) \xrightarrow{n \rightarrow \infty} \gamma.$$

Замечание. Этот доверительный интервал не очень хорош, потому что может включать значения меньше 0.

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 обязательно принадлежит $[0, 1]$.

Задание: Построить хорошие доверительные интервалы в моделях Бернулли и Пуассона, решив неравенства с квадратными уравнениями.

Задание: Построить доверительный интервал для параметра λ экспоненциального распределения $\text{Exp}(\lambda)$.

3.6. Использование SE для проверки гипотез и построения доверительных интервалов

Пусть оценка $\hat{\theta}_n$ имеет (асимптотически) нормальное распределение и является асимптотически несмещенной. Тогда доверительный интервал уровня γ для θ (т.е. такой интервал, в котором лежит γ всех значений величины) задается как

$$\hat{\theta}_n \pm c_\gamma \sqrt{D\hat{\theta}_n},$$

где c_γ — $(\gamma + 1)/2$ -квантиль стандартного нормального распределения. К примеру, для $N(0, 1)$ и 95%-квантили это был бы интервал $(-1.96; 1.96)$, а так нужно передвинуть его на среднее и растянуть на корень из дисперсии.

Но стандартное отклонение $\sqrt{D\hat{\theta}_n}$ распределения $\hat{\theta}_n$ можно оценить как SE (standard error, стандартная ошибка). Значит, доверительный интервал будет иметь вид

$$\hat{\theta}_n \pm c_\gamma SE.$$

Аналогично, статистика критерия $H_0 : \theta = \theta_0$ будет иметь вид

$$t = (\hat{\theta}_n - \theta_0) / SE(\hat{\theta}_n),$$

которая имеет асимптотически нормальное распределение $N(0, 1)$ с ‘идеальным’ значением в нуле.

Заметим, что SE играет роль ‘сигмы’ распределения оценки.

4. Критерии проверки гипотезы о согласии с видом распределения

4.1. Критерий χ^2

По выборке возможно проверить гипотезу о виде распределения случайной величины, реализацией которой является выборка. Для проверки гипотезы согласия с видом произвольного *дискретного* распределения используется асимптотический критерий χ^2 («chi-squared test for goodness of fit»).

4.1.1. 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 (следовательно, критическая область только справа).

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

Определение. Критерий применим, если $\alpha_I = \alpha$ или $\alpha_I \approx \alpha$ с достаточной степенью точности.

Замечание. Поскольку критерий асимптотический, с достаточной (тому, кто дает такие рекомендации) 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.

Замечание. Почему бы не подстраховаться и не объединить состояния так, чтобы было > 10 ?
Ответ: теряем в мощности.

Задание Привести пример, демонстрирующий потерю мощности.

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

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

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. Если нужно подстраховаться от подгонки (искусственно составленных под гипотезу выборок), то критическую область можно выбрать с двух сторон, слева и справа. Например, чтобы отверглась гипотеза о $p = 0.5$ для альтернирующей (и явно не случайной) последовательности $\mathbf{x} = (0, 1, 0, 1, \dots)$ имеет $T = 0$. Однако, если мы не подозреваем данные в обмане, то так не делают.

4.1.2. Распределение с неизвестными параметрами

В случае сложной гипотезы $\mathcal{P} \in \{\mathcal{P}(\boldsymbol{\theta})\}_{\boldsymbol{\theta} \in \Theta}$, $\boldsymbol{\theta} = (\theta_1, \dots, \theta_r)^\top$, следует найти оценку $\hat{\boldsymbol{\theta}}_{\text{MLE}}$ 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,

$$T = \sum_{i=1}^k \frac{(\nu_i - np_i(\hat{\boldsymbol{\theta}}_{\text{MLE}}))^2}{np_i} \xrightarrow{d} \chi^2(k - r - 1).$$

Важно: параметр можно считать известным, только если его значение выбрано без знания, какая получилась выборка.

Оценки по методу минимума хи-квадрат Предельное распределение статистики критерия не поменяется, если вместо оценки максимального правдоподобия подставить любую другую оценку с тем же предельным распределением. Рассмотрим оценки по минимуму хи-квадрат:

$$\boldsymbol{\theta}_{\text{minChiSq}} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^k \frac{(\nu_i - np_i(\boldsymbol{\theta}))^2}{np_i} \xrightarrow{d} \chi^2(k - r - 1).$$

Утверждение (без доказательства) заключается в том, что в условиях регулярности оценки, полученные по методу максимального правдоподобия, эквивалентны оценкам, полученным по методу минимума хи-квадрат. Таким образом, в критерии хи-квадрат можно использовать статистику в виде

$$T = \min_{\boldsymbol{\theta}} \sum_{i=1}^k \frac{(\nu_i - np_i(\boldsymbol{\theta}))^2}{np_i} \xrightarrow{d} \chi^2(k - r - 1).$$

5. Критерий Колмогорова-Смирнова согласия с видом распределения

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 точный 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. Распределение D_n , оно разное для разных n , но не зависит от распределения.

Распределение не зависит от cdf_0 , так как любое распределение можно привести, например, к равномерному, монотонным преобразованием: для любой cdf_0 верно $\text{cdf}_0^{-1}(\xi) \sim U(0, 1)$, как будто, мы проверяем гипотезу $H_0 : \text{cdf}_0^{-1}(\xi) \sim U(0, 1)$.

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

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

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

6. Визуальное определение согласия с распределением

6.1. 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)  
}
```

6.2. 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 $y = \sqrt{\widehat{D\xi}} \cdot x + \widehat{E\xi}$.

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/>

Часть IV.

Корреляционный анализ

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

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

1. Вероятностная независимость

1.1. Визуальное определение независимости

- Поскольку при $p_\eta(y_0) \neq 0$

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

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 строках / 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:

$$\begin{array}{ccc} 1 & 3 & 2 \\ 2 & 5 & 3 \\ 9 & 20 & 11 \end{array}$$

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 $\hat{p}_{\cdot j} = \nu_{\cdot j}/n$.

1. Вероятностная независимость

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. Сводится к гипотезе о независимости, так как независимость эквивалентна равенству условных распределений. \square

Замечание. На маленьких выборках ($n < 50$ или $np_{ij} < 5$) 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 перестановкой¹ 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. Однако есть требование дискретности (конечного числа значений).

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

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

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

или коэффициент сопряженности Пирсона

$$\hat{p}^2 := \frac{\chi^2}{\chi^2 + n} = \frac{\hat{r}^2}{\hat{r}^2 + 1}.$$

(тогда 1 никогда не достигается).

Заметим, что $\hat{r}^2 := \frac{\chi^2}{n}$ является оценкой следующей меры зависимости (меры сопряженности) для двумерного дискретного распределения, задаваемого набором p_{ij} :

$$r^2 = \sum_{i=1}^K \sum_{j=1}^S \frac{(p_{ij} - p_{i\cdot}p_{\cdot j})^2}{p_{i\cdot}p_{\cdot j}}$$

¹[https://en.wikipedia.org/wiki/Resampling_\(statistics\)#Permutation_tests](https://en.wikipedia.org/wiki/Resampling_(statistics)#Permutation_tests)

2. Линейная / нелинейная зависимость

2.1. Определение вида зависимости

Пусть теперь ξ, η — количественные признаки.

Напомним, условное математическое ожидание $E(\eta \mid \xi)$ является такой функцией от ξ , на которой достигается минимум $\min_{\hat{\eta} \in \{\varphi(\xi)\}} E(\eta - \hat{\eta})^2$.

Определение. Определим функцию условного математического ожидания

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

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

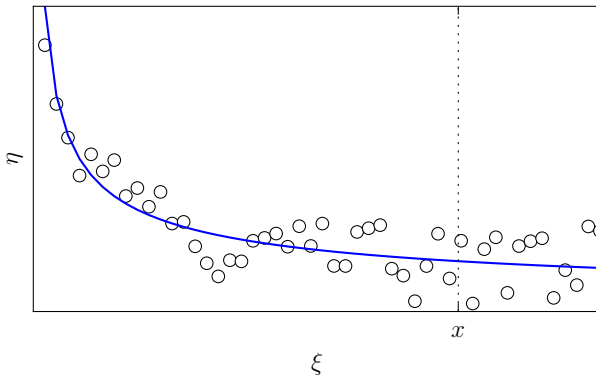


Рис. 2.1.: Нелинейная зависимость

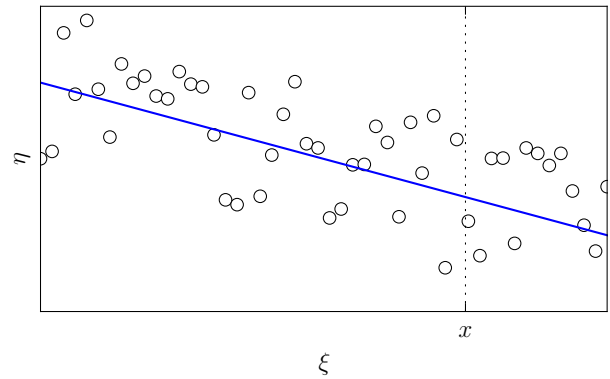


Рис. 2.2.: Линейная зависимость

2.2. Коэффициент корреляции Пирсона

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).

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

Предложение. Для линейно зависимых данных, конечно, $\rho = \text{sign } b$.

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

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

□

Предложение.

$$\rho^2(\xi, \eta) = 1 - \frac{\min_{\hat{\eta} \in \{a+b\xi\}} \mathbb{E}(\eta - \hat{\eta})^2}{D\eta}.$$

2.2.1. Оценка коэффициента корреляции

Оценка коэффициента корреляции строится стандартным методом подстановки в формулу для корреляции двумерного эмпирического распределения, в котором каждая пара значений $(x_i, y_i)^T$, $i = 1 \dots, n$, имеет вероятность $1/n$.

В знаменателе стоят дисперсии, поэтому они в оценке просто заменяются на выборочные дисперсии.

Для оценки ковариации $\text{cov}(\xi, \eta) = \mathbb{E}(\xi - \mathbb{E}\xi)(\eta - \mathbb{E}\eta) = \mathbb{E}\xi\eta - \mathbb{E}\xi\mathbb{E}\eta$ можно использовать два варианта (одной и той же) оценки, поэтому получим:

$$\hat{\rho}(\xi, \eta) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y} = \frac{\sum_{i=1}^n x_i y_i - \bar{x}\bar{y}}{s_x s_y}.$$

2.2.2. Значимость коэффициента корреляции

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

Пусть $(\xi, \eta)^T \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Тогда при $H_0 : \rho = 0$ статистика критерия имеет вид и распределение:

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

Идеальное значение — 0, критическая область двухсторонняя. Если предположения о нормальности $(\xi, \eta)^T$ нет, а гипотеза не о некоррелированности, а о независимости, то критерий становится асимптотическим (т.е. им все равно можно пользоваться).

Проверим теперь гипотезу $H_0 : \rho = \rho_0$ (чаще проверяют $H_0 : \rho > \rho_0$). Тогда применяется z -преобразование Фишера

$$z = \frac{1}{2} \ln \frac{1+\rho}{1-\rho}, \quad z_0 = \frac{1}{2} \ln \frac{1+\rho_0}{1-\rho_0}.$$

Если $(\xi, \eta)^T \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$,

$$T = \sqrt{n-3}(z - z_0) \xrightarrow{d} N(0, 1).$$

2.3. Происхождение и сравнение других мер зависимости

Если K — линейное пространство, то теорема Пифагора принимает вид

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

где $\hat{\eta}^* = \text{argmin}_{\hat{\eta} \in K} \mathbb{E}(\eta - \hat{\eta})^2$.

Откуда можно записать меру аппроксимации линейным пространством K как

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

Если $K = \mathcal{L} = \{a\xi + b\}$, то полученная величина является квадратом коэффициентом корреляции ρ^2 :

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

ρ — коэффициент корреляции Пирсона.

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

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

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*:

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

Сравнивая полученные формулы, мы получаем, что коэффициент корреляции измеряет, насколько хорошо случайную величину η можно приблизить линейной функцией от ξ , а корреляционное отношение — произвольной (измеримой) функцией от ξ .

2.3.1. Свойства корреляционного отношения

1. $r_{\eta|\xi}^2 \in [0, 1]$.
2. $\eta \perp\!\!\!\perp \xi \implies r_{\eta|\xi}^2 = 0$.
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. $(\xi, \eta)^T \sim N(\mu, \Sigma) \implies r_{\eta|\xi}^2 = \rho^2(\eta, \xi)$.

2.3.2. Выборочное корреляционное отношение

По разложению дисперсии,

$$D\eta = \mathbb{E}(\eta - \mathbb{E}\eta)^2 = \underbrace{\mathbb{E}(\mathbb{E}(\eta | \xi) - \mathbb{E}\eta)^2}_{D\mathbb{E}(\eta|\xi)} + \mathbb{E}(\eta - \mathbb{E}(\eta | \xi))^2.$$

Перейдем на выборочный язык. Пусть дана выборка

$$\begin{pmatrix} x_1 \\ y_1 \end{pmatrix}, \dots, \begin{pmatrix} x_n \\ y_n \end{pmatrix}.$$

Сгруппируем её:

$$\begin{array}{c|ccc} x_1^* & y_{11} & \dots & y_{1n_1} \\ \vdots & \vdots & \ddots & \vdots \\ x_k^* & y_{k1} & \dots & y_{kn_k} \end{array}$$

Пусть ξ — дискретная случайная величина со значениями (x_1^*, \dots, x_k^*) . Тогда, учитывая

$$\bar{y}_i = \bar{y}_{i\cdot} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} = \hat{\mathbb{E}}(\eta | \xi = x_i^*),$$

на выборочном языке получаем (домножив на n):

$$\underbrace{\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2}_{\text{total sum of squares}} = \underbrace{\sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2}_{\text{межгрупповой разброс}} + \underbrace{\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}_{\text{внутригрупповой разброс}}$$

$$ns_y^2 = ns_{y|x}^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2.$$

Отсюда, так как, $r_{\eta|\xi}^2 = \text{DE}(\eta \mid \xi) / \text{D}\eta$,

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

3. 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.

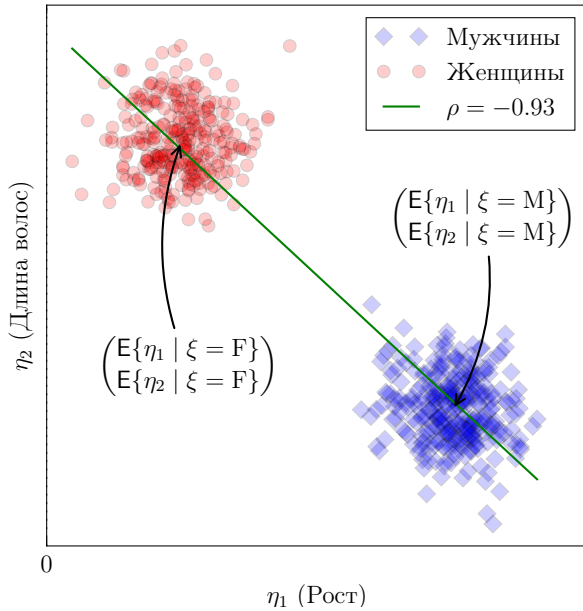


Рис. 3.1.: Исходные данные (бимодальность)

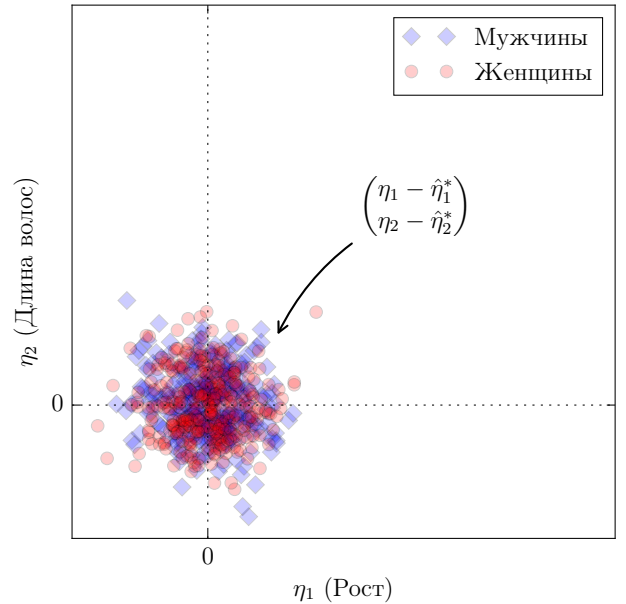


Рис. 3.2.: Центрированные данные

Пример. Возможна и ситуация как на (3.3), где определено $\rho(\eta_1, \eta_2) > 0$, но $\rho(\eta_1, \eta_2 \mid \xi) < 0$.

3. Частная корреляция

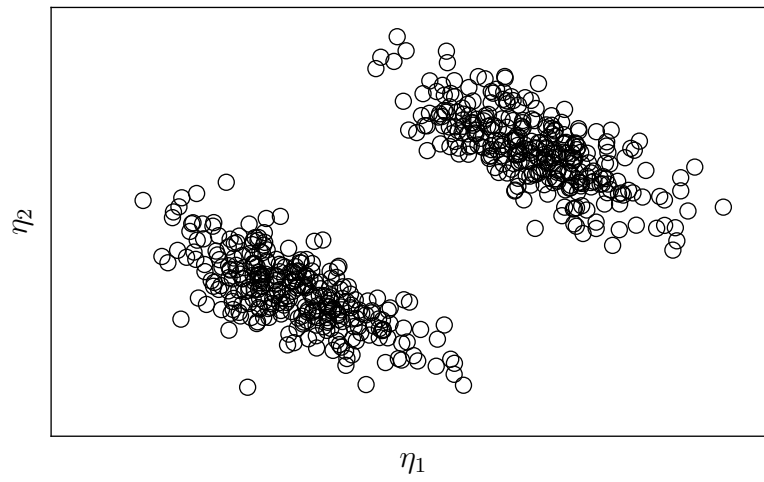


Рис. 3.3.: $\rho(\eta_1, \eta_2) > 0$, но $\rho(\eta_1, \eta_2 \mid \xi) < 0$

4. Зависимость между порядковыми признаками

Пусть признаки порядковые, т.е. их значения можно только сравнивать. Это означает, что нельзя читать, например, математическое ожидание, плотность и пр., но понятие функции распределения, основанное на сравнении, корректно.

Оценки значений функции распределения строятся по выборке на рангах случайных величин. Поэтому в случае порядковых признаков рассматриваются оценки, основанные на рангах. В частности, ранговые коэффициенты корреляции. Заметим, что ранговые характеристики хорошо работают на выборках без совпадающих наблюдений.

4.1. Ранговый коэффициент Спирмана

Определение. Ранговый коэффициент Спирмана есть

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

Замечание. Для непрерывной функции распределения, $\text{cdf}_\xi(\xi) \sim U(0, 1)$, 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

$$\text{cdf}_{\xi_n}(x_i + 0) = \frac{\text{rk } x_i}{n}, \quad \text{cdf}_{\eta_n}(y_i + 0) = \frac{\text{rk } y_i}{n} = \frac{T_i}{n}.$$

Тогда будет справедливо следующее

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

$$\hat{\rho}_S = \frac{1/n \cdot \sum_{i=1}^n R_i T_i - \bar{R} \bar{T}}{\sqrt{1/n \cdot \sum_{i=1}^n (R_i - \bar{R})^2} \sqrt{1/n \cdot \sum_{i=1}^n (T_i - \bar{T})^2}}.$$

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 $|\rho_S|$, тем более явно выражена зависимость y от x в виде некоторой монотонной функции.

4.1.1. Согласованность ρ и ρ_S

Для количественных признаков, мера монотонной зависимости ρ_S не согласована с мерой линейной зависимости ρ в том же смысле, что ρ и мера функциональной зависимости $r_{\xi|\eta}$, где, в частности, $\rho \leq r_{\xi|\eta}$, а равенство достигается в случае линейной зависимости (линейного условного математического ожидания).

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

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

- С точностью до погрешности, по значению, ρ и ρ_S — это одно и то же (см. 4.1)

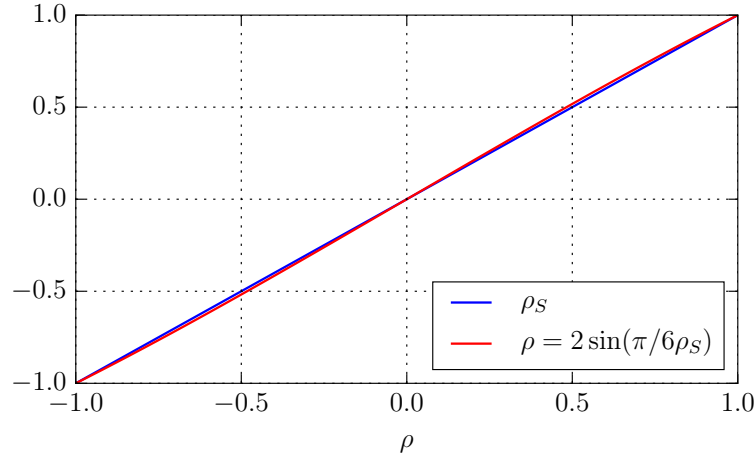


Рис. 4.1.: $\rho \approx \rho_S$

Значит, в случае нормального распределения можем сравнить оценки $\hat{\rho}$ и $\hat{\rho}_S$ между собой.

- Выборочную дисперсию оценок сравнить довольно сложно. Тем не менее, можем заметить, что $\hat{\rho}_S$ более устойчив к аутлаерам (см. 4.2). 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.

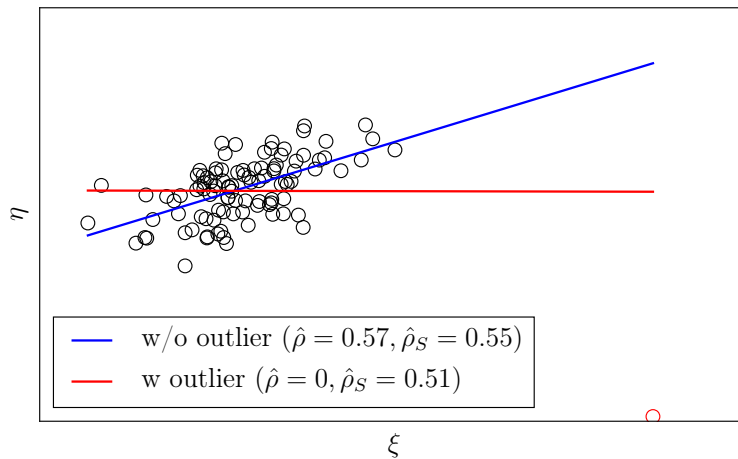


Рис. 4.2.: $\hat{\rho}$ до и после добавления аутлаера

- Монотонным преобразованием можем всегда сделать так, чтобы ρ изменился (например, возведя в квадрат); при монотонном преобразовании, однако, не меняется ρ_S (см. 4.3).

Значит, чтобы узнать ρ исходных (нормальных) данных, можно не выполнять обратного преобразования, а сразу посчитать ρ_S .

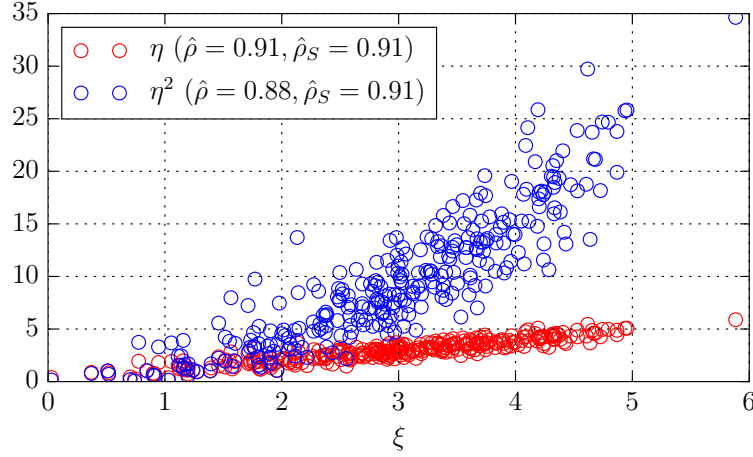


Рис. 4.3.: Монотонное преобразование нормальных данных

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

Определение. Пусть $(\xi_1, \eta_1)^\top \perp (\xi_2, \eta_2)^\top \sim \mathcal{P}_{\xi, \eta} \sim (\xi, \eta)^\top$; 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 $(\xi, \eta)^\top \sim N(\boldsymbol{\mu}, \Sigma)$, 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).

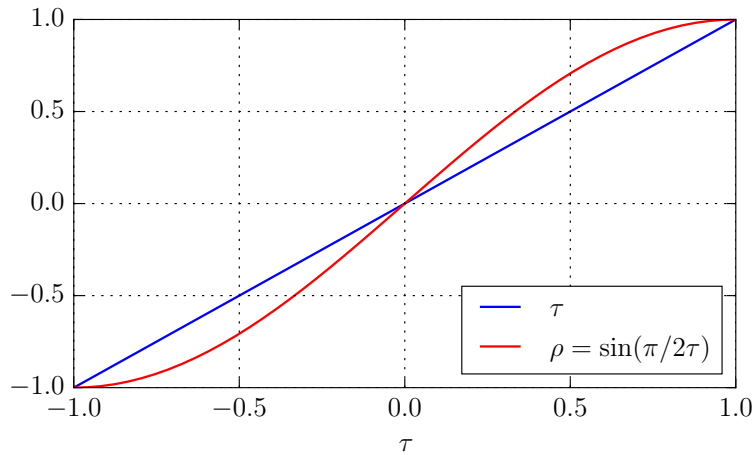


Рис. 4.4.: ρ и τ

4. Зависимость между порядковыми признаками

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.