

СОВРЕМЕННЫЕ МЕТОДЫ ОБРАБОТКИ ДАННЫХ

Грибанов К.Г. для группы МЕНМ-190501

Случайные величины и их распределения

Теория вероятностей и математическая статистика составляют фундамент знаний, необходимых для понимания природы любых данных, с которыми приходится иметь дело специалистам из области естественных наук.

Поскольку события реального мира не являются числами, то для того чтобы связывать события и числа, чтобы использовать математический аппарат, необходимо в общих чертах иметь представление о теории множеств. Практически, мы все всегда манипулируем понятиями теории множеств, когда приписываем событию или объекту реального мира какое-нибудь число, целое или действительное. Простой пример: алфавитный список группы студентов, посещающих лекции по данному предмету ставит в соответствие реальному человеку порядковый номер в списке, т. е. целое число. Если всех студентов взвесили, то поставили в соответствие действительное число, округленное до точности измерения, т. е. вес, например в килограммах и десятых килограмма. Если в первом случае каждому студенту поставлено в соответствие уникальное число, т. е. номер, то в случае взвешивания может оказаться, что у кого-то вес совпадет и число не будет уникальным. Здесь неявно приведено два примера разных отображений из теории множеств.

Случайной величиной называется величина (число), которая в результате одного наблюдения (опыта) может принимать случайным образом только одно значение, при этом условия проведения опыта фиксированы для серии наблюдений (опытов). Случайную величину в данном курсе будем обозначать, как это делается во множестве учебников, прописной латинской буквой, например X , а её значения (называющиеся *реализациями*), получаемые при наблюдении строчной буквой x . СВ может быть дискретной, непрерывной, скалярной или векторной. *Вероятность* это мера того, как часто та или иная СВ наблюдается в серии экспериментов, это тоже число, но действительное и обязательно лежащее в интервале $[0; 1]$. *Дискретная случайная величина* может принимать значения только из некоторого конечного или счетного множества. *Непрерывная случайная величина* принимает любые значения на некотором плотном отрезке или области, например $X \in [A; B]$, причем отрезок может быть и бесконечным.

Распределением (законом распределения) дискретной СВ является совокупность вероятностей всех ее возможных значений, т. е. $P(x_\nu)$, $\nu = 1, 2, \dots, N$. Примером такого распределения является биномиальное распределение, которое описывает число m выпадений герба в n бросках, если форма монеты дает вероятность выпадения герба q . Кроме того, данное распределение может описывать число отказов в испытании серии новых изделий. Формула биномиального распределения следующая:

$$P(m) = \frac{n!}{m!(n-m)!} q^m (1-q)^{n-m}, \quad m=0, 1, 2, \dots, n$$

На Рис.1 приведен пример биномиального распределения для $n=20$ и $q=0.3$. (Как создать такую картинку средствами языка сценариев R изложено в отдельной презентации «Распределения дискретных и непрерывных случайных величин и знакомство с ними в R»)

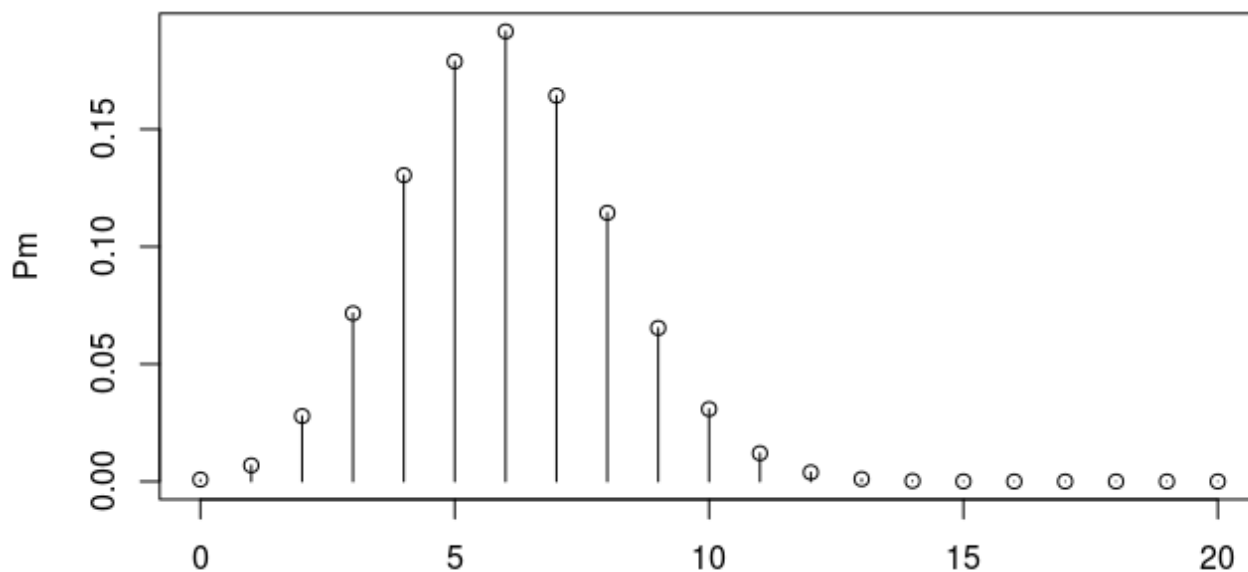


Рис. 1. Пример биномиального распределения для $n=20$ и $q=0.3$.

Вероятности дискретного распределения подчиняются условию нормировки:

$$\sum_{m=0}^n P(m) = 1 .$$

Распределение дискретной СВ можно также представить в виде ступенчатой функции распределения

$$F(m) = P\{M < m\} = \sum_{m' < m} P(m') ,$$

график которой приведен для того же биномиального распределения на Рис.2. Следует понимать условность такой функции и ее ступенчатого графика для дискретного распределения, поскольку значения аргумента, откладываемые по оси абсцисс на Рис.2 могут принимать только дискретные значения, а значений абсциссы между ними, на «полках» ступенчатой функции нет во множестве случайной величины.

Распределение (закон распределения) непрерывной СВ характеризуют плотностью распределения вероятности $f(x)$ или функцией распределения $F(x)$. Плотность [распределения] вероятности скалярной СВ, является производной от функции распределения, т. е.

$$f(x) = (F(x))' .$$

Плотность вероятности позволяет вычислить вероятность попадания СВ в заданный интервал, т. е.

$$P\{X \in [x - \delta; x + \delta]\} = \int_{x - \delta}^{x + \delta} f(x') dx' .$$

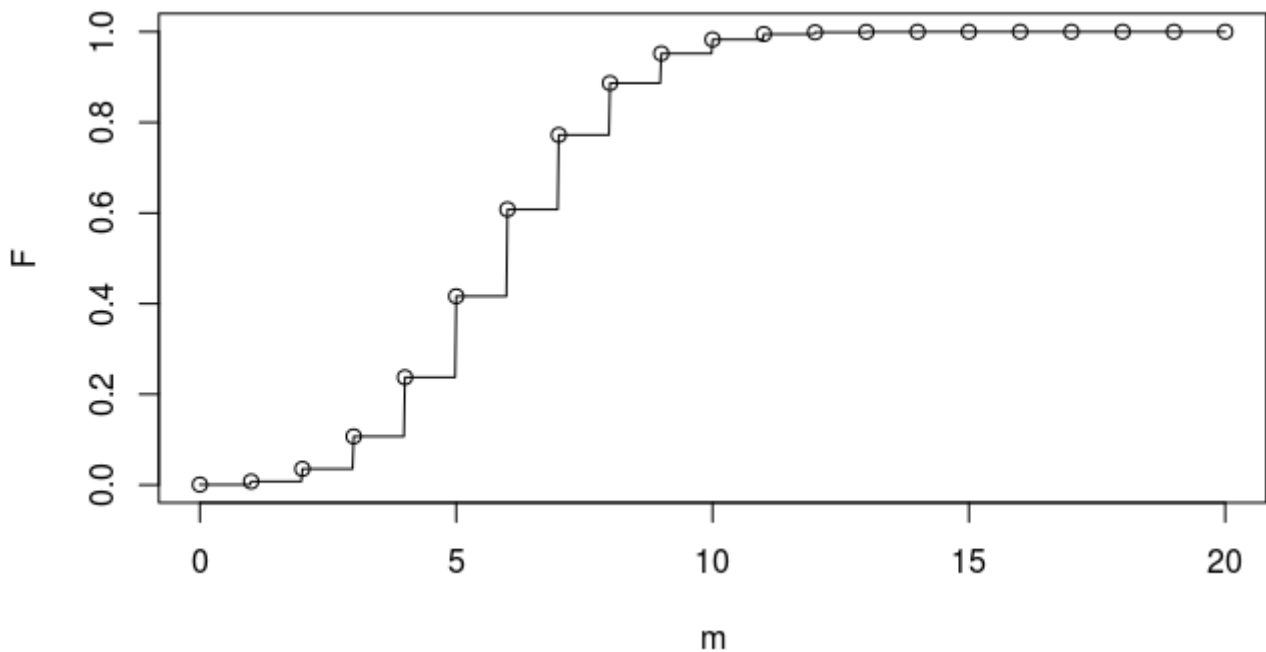


Рис. 2. Ступенчатая функция распределения для биномиального распределения, показанного на Рис. 2

Функция же распределения представляет собой вероятность того, что СВ $X \in [A; B]$ не превысит некоторое значение x , т.е.

$$F(x) = P\{X \leq x\} = \int_A^x f(x') dx' ,$$

здесь штрихованная переменная под интегралом используется для того, чтобы отличать ее от верхнего предела интегрирования. Вот некоторые свойства функций плотности вероятности и распределения:

1. $f(x) \geq 0, \forall x \in [A; B]$, т. е. $f(x)$ неотрицательная;

$$2. \int_A^B f(x) dx = 1, \text{ условие нормировки;}$$

$$3. F(A) = 0, \quad F(B) = 1;$$

4. Если $x_1 \geq x_2$, то $F(x_1) \geq F(x_2)$, т. е. функция $F(x)$ монотонно возрастает.

На Рис. 3 в качестве примера приведены $f(x)$ и $F(x)$ для гамма-распределения (с которым подробнее следует познакомиться с помощью презентации «Распределения дискретных и непрерывных случайных величин и знакомство с ними в R»).

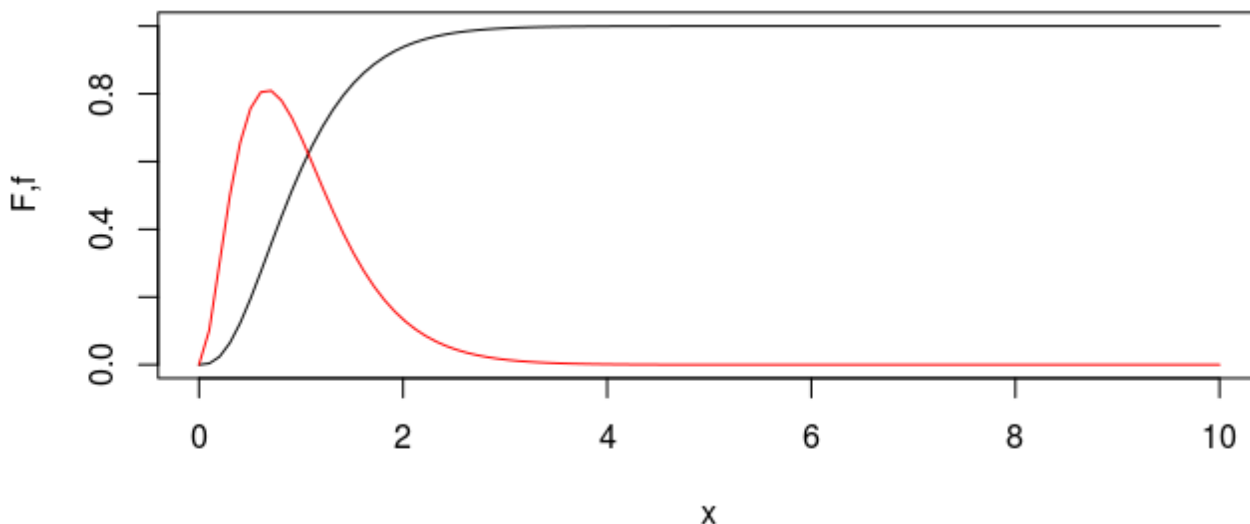


Рис. 3. Пример плотности вероятности для гамма-распределения (красная кривая) и функции распределения (черная кривая).

Векторная случайная величина это совокупность взаимосвязанных скалярных СВ, $X = (X_1, X_2, \dots, X_n)^T$ (здесь T — это знак транспонирования, чтобы представлять вектор в виде матрицы столбца).

Моменты распределения случайных величин

Для решения практических задач часто достаточно только двух первых моментов характеризующих положение — *математическое ожидание* и рассеяние — *дисперсия*. *Математическое ожидание* для дискретной СВ

$$M_X = \sum_{n=1}^N x_n P(x_n),$$

для непрерывной СВ

$$M_X = \int_A^B x f(x) dx .$$

Дисперсия СВ представляет собой второй центральный момент распределения и определяется следующим образом:

для дискретной СВ
$$D_X = \sum_{n=1}^N (x_n - M_X)^2 P(x_n) ,$$

для непрерывной СВ
$$D_X = \int_A^B (x - M_X)^2 f(x) dx .$$

Размерность дисперсии не совпадает с размерностью СВ, поэтому на практике используется *среднеквадратическое отклонение* (СКО) СВ:

$$\sigma_X = \sqrt{D_X} .$$

С учетом этой формулы, дисперсию часто обозначают как σ_X^2 . На практике часто используют «правило трех сигм», имея в виду, что все реализации СВ почти наверное попадают в интервал $[M_X - 3\sigma_X; M_X + 3\sigma_X]$.

Распределение векторной СВ $X = (X_1, X_2, \dots, X_n)^T$ характеризуется векторным МО и ковариационной матрицей K_X . Математическое ожидание векторной СВ представляет собой вектор состоящий из математических ожиданий компонентов вектора, т.е.

$$M_X = (M_{X_1}, M_{X_2}, \dots, M_{X_n})^T .$$

Ковариационная матрица векторной СВ состоит из дисперсий компонентов на главной диагонали и моментов ковариации пар компонентов в качестве внедиагональных компонентов:

$$K_X = \begin{pmatrix} D_{X_1} & K_{X_1, X_2} & \dots & K_{X_1, X_n} \\ K_{X_2, X_1} & D_{X_2} & \dots & K_{X_2, X_n} \\ \dots & \dots & \dots & \dots \\ K_{X_n, X_1} & K_{X_n, X_2} & \dots & D_{X_n, X_n} \end{pmatrix} .$$

Коэффициенты ковариации можно определить в терминах математического ожидания от произведения центрированных пар компонентов векторной СВ, т.е.

$$K_{X_i, X_j} = E[(x_i - M_{X_i})(x_j - M_{X_j})] ,$$

где E это оператор математического ожидания, если $i=j$, то формула переходит в формулу для дисперсии на главной диагонали.

Простой пример. *Равномерное распределение* на отрезке $[A; B]$ имеет плотность вероятности

$$f(x) = \begin{cases} \frac{1}{B-A}, & x \in [A; B] \\ 0, & x \notin [A; B] \end{cases},$$

и функцию распределения

$$F(x) = \begin{cases} 0, & x < A \\ \frac{x-A}{B-A}, & x \in [A; B] \\ 1, & x > B \end{cases}.$$

На Рис. 4 показаны обе функции для $A=1, B=6$.

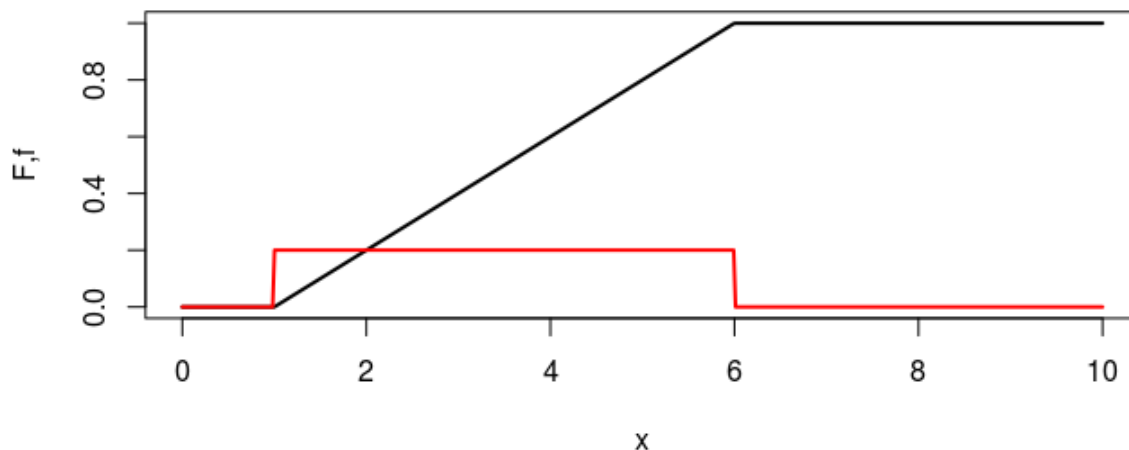


Рис. 4. Функция распределения (черная) и плотность вероятности (красная) для равномерного распределения на отрезке $[1; 6]$.

Предлагается самостоятельно получить выражения для M_x и D_x для равномерного распределения на отрезке $[A; B]$, а также воспроизвести Рис.4 средствами R.

Задачи математической статистики

Математическая статистика связана с теорией вероятностей и в ней ставятся и решаются задачи изучения закономерностей и получения выводов о СВ по выборке данных наблюдений, принадлежащих совокупности СВ, называемой *генеральной*.

Генеральная совокупность — совокупность объема N всех физических объектов одного вида или значений СВ с определенным распределением. Число N может быть конечно, а может и $N \rightarrow \infty$.

Выборка — совокупность элементов размера $n < N$, случайно полученная из ГС при некоторых условиях эксперимента. Выборка может быть *повторной* — отобранные элементы возвращаются в ГС, или *бесповторной* — элементы не возвращаются в ГС. Задача математической статистики по выборке получить сведения о ГС. ГС может быть *стабильной* или *динамической*, т. е. изменяющейся со временем. Чтобы не углубляться в теорию случайных процессов, в данном курсе рассматриваются только стабильные ГС. Выборка должна быть *репрезентативной (представительной)*, что реализуется если она реализуется случайно, т. е. все элементы ГС имеют одинаковые шансы попасть в выборку. Выборка должна быть *достаточной*, т. е. иметь достаточный объем, чтобы обеспечить хорошую точность сведений о ГС. Это не всегда удается достичь и приходится довольствоваться выборкой малого объема.

Из выборки можно получить *статистическое распределение* СВ. Для дискретной СВ $x_q, q=1,2,3,\dots,Q$ при наличии выборки объема n , статистическое распределение определяется как

$$p_q = \frac{n_q}{n}, \quad q=1,2,\dots,Q,$$

где n_q это число наблюдений величины x_q в выборке объема n . Естественно выполнение следующих соотношений:

$$\sum_n n_q = n, \quad \sum_n p_q = 1.$$

Статистическое распределение дискретной СВ можно отобразить на графике в стиле Рис.1.

Для непрерывной СВ, при наличии выборки $x_i, i=1,2,3,\dots,n$ можно построить *гистограмму распределения* или *эмпирическую функцию распределения*. Для этого диапазон реализации $\Delta = x_{max} - x_{min} + \varepsilon$ делится на K отрезков одинаковой длины $h = \Delta/K$, K рекомендуется выбирать в интервале от 6 до 20, а в случае большой выборки $K \approx n^{1/3}$. Значение $\varepsilon < h$ выбирается из соображений удобства представления гистограммы. Выбираются середины отрезков $x_k, k=1,2,\dots,K$, вычисляются n_k - числа попадания СВ из выборки в интервал $[x_k - h/2; x_k + h/2]$, а затем вероятности попадания в этот интервал

$$P_k = \frac{n_k}{n}, \quad k=1,2,3,\dots,K.$$

Можно построить гистограмму из прямоугольников высотой P_k и шириной h , однако если нормировать высоту прямоугольников $f_k = P_k/h$, можно

построить статистическую плотность вероятности, площадь под гистограммой в этом случае будет удовлетворять условию нормировки

$$\sum_{k=1}^K f_k h = 1 .$$

Для построения эмпирической функции распределения на основе выборки строится вариационный ряд, т. е. реализации x_i записываются в возрастающем порядке, при этом получается новая последовательность $x_q, q=1,2,\dots,n$, которая используется для вычисления ФР:

$$F(x_q) = P\{x < x_q\} = \frac{(q-1)}{n}, \quad q=1,2,\dots,n .$$

При $n \rightarrow \infty$, $f_k \rightarrow f(x_k)$, $F(x_q) \rightarrow F(x)$. Примеры обоих эмпирических функций для нормального распределения показаны на Рис.5.

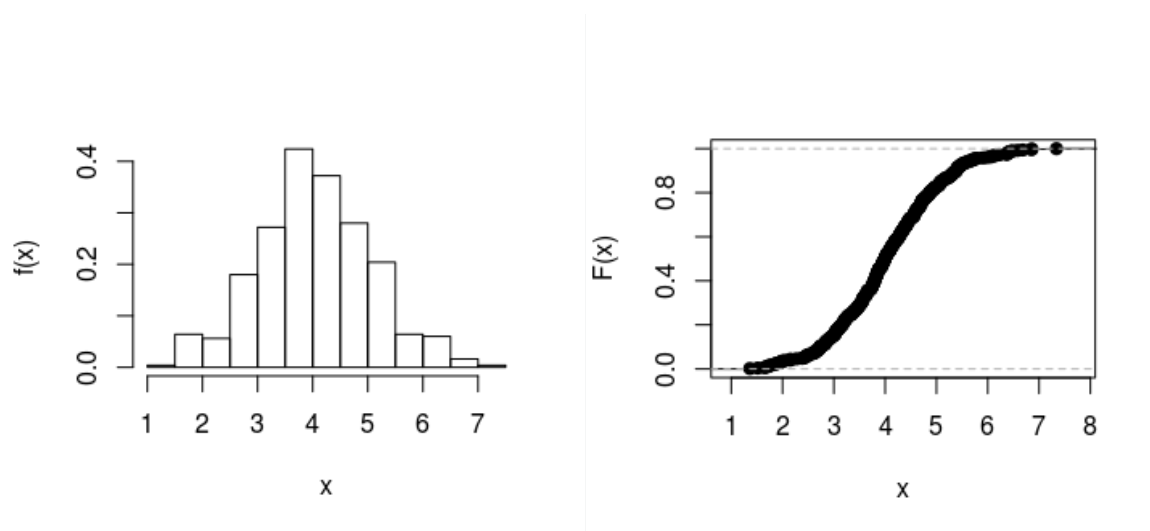


Рис.5. Примеры статистической плотности вероятности и функции распределения. Генерация выборки нормально распределенной СВ выполнена с помощью функции `rnorm()`, а построение графиков выполнено в R с помощью функций `hist()` и `ecdf()`.

Точечное оценивание параметров распределения случайных величин

Оцениваемые на основе ограниченной выборки параметры распределения случайной величины являются неточными и называются *статистическими оценками*, которые зависят от выборки и сами являются случайными величинами. *Точечная оценка* это величина, зависящая от выборки и выражаемая одним числом или вектором. Пусть СВ имеет распределение с параметром θ , *точечную оценку* или просто *оценку* этого параметра будем

обозначать $\hat{\theta}$. Для получение точечных оценок используются три основных метода: метод максимального правдоподобия (ММП), метод наименьших квадратов (МНК) и метод моментов (ММ). Поскольку оценка это СВ, то и она характеризуется математическим ожиданием

$$M_{\hat{\theta}} = E [\hat{\theta}]$$

и дисперсией $\sigma_{\hat{\theta}}^2$. Иногда удается построить плотность вероятности оценок параметра $f(\theta)$ и если такая функция построена корректно, то можно определить *несмещенную точечную оценку параметра θ* и её дисперсию как

$$\hat{\theta} = \int_{-\infty}^{\infty} \theta f(\theta) d\theta ,$$

$$\sigma_{\hat{\theta}}^2 = \int_{-\infty}^{\infty} (\theta - \hat{\theta})^2 f(\theta) d\theta .$$

К точечной оценке предъявляются требования по её *несмещённости, эффективности, состоятельности* и *достаточности*.

Несмещенность оценки означает что математическое ожидание оценки равно точному параметру, т.е.

$$E [\hat{\theta}] = \theta .$$

Эффективность оценки означает, что ее дисперсия является минимальной в классе всех несмещенных оценок параметра θ :

$$\sigma_{\hat{\theta}}^2 = \min .$$

Состоятельная оценка сходится по вероятности к истинному значению при увеличении объема выборки, т.е.

$$P \{ |\hat{\theta} - \theta| < \delta \} \xrightarrow[n \rightarrow \infty, \delta \rightarrow 0]{} 1 .$$

Достаточной называется оценка, которая содержит такое же количество информации о неизвестном параметре θ , что и элементы выборки.

Получение информации с помощью $f(\theta)$ затруднено тем, что такую функцию редко удается построить. Свойства оценок скалярного параметра, рассмотренные выше, можно распространить и на векторный параметр θ .

Интервальное оценивание

Интервальная оценка это статистическая оценка в виде совокупности точек или множества. Обычно это *доверительный интервал* $[\theta_H; \theta_B]$, такой что

$$P\{\theta_H \leq \theta \leq \theta_B\} = \gamma,$$

где γ это *доверительная вероятность*, θ_H и θ_B - *нижняя и верхняя доверительные границы* соответственно. Нижняя или верхняя границы могут оказаться и бесконечными. Точность интервальной оценки характеризуется размахом доверительного интервала, а ее достоверность или надежность — доверительной вероятностью.

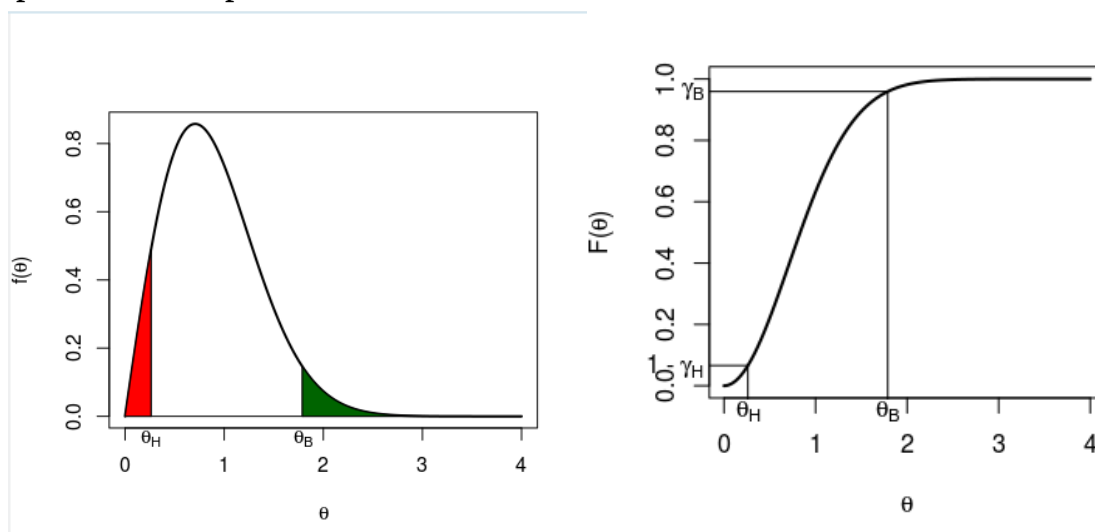


Рис. 6. Определение доверительных границ для параметра θ по плотности вероятности (слева) или функции распределения (справа). На графике слева площади закрасенные красным, белым и зелёным равны вероятностям $1-\gamma_H$, γ и $1-\gamma_B$ соответственно.

Оценки доверительных границ двухстороннего доверительного интервала определяются следующими зависимостями:

$$P\{\theta \geq \theta_H\} = \gamma_H, \quad P\{\theta \leq \theta_B\} = \gamma_B,$$

$$\int_{-\infty}^{\theta_H} f(\theta) d\theta = 1-\gamma_H, \quad \int_{\theta_B}^{\infty} f(\theta) d\theta = 1-\gamma_B,$$

$$F(\theta_H) = 1-\gamma_H, \quad F(\theta_B) = \gamma_B.$$

Из правого графика на Рис. 6 можно получить условие:

$$\gamma = \gamma_B - (1-\gamma_H) \quad \text{или} \quad \gamma_H + \gamma_B = 1 + \gamma.$$

При симметричной $f(\theta)$ для упрощения равенств выше величины y_H и y_B можно выбирать равными:

$$y_H = y_B = 0.5(1 + \gamma) .$$

Не всегда удастся построить плотность вероятности или функцию распределения для оценки параметра, поэтому для интервального оценивания некоторых параметров распределения СВ, например МО и дисперсии нормального распределения, удастся подобрать специальные функции (*статистики*) g , которые зависят от результатов наблюдения $x_i, i=1,2,\dots,n$ и от точечной оценки параметра θ : $g = \psi(\{x_i\}, \theta)$, для которых удастся построить плотность вероятности $\varphi(g)$. В этом случае для функции g определяются доверительные границы $[G_H; G_B]$, а оценки доверительных границ $[\theta_H; \theta_B]$ уже для параметра θ вычисляются как решения уравнений

$$\psi(\{x_i\}, \theta_H) = G_H , \quad \psi(\{x_i\}, \theta_B) = G_B .$$

Для определения доверительных границ θ_H и θ_B некоторых параметров распределения СВ используются заранее подготовленные таблицы квантилей распределения функции g в зависимости от значений n и γ . Для интервального оценивания математического ожидания и дисперсии нормальной СВ используются табличные значения квантилей распределений Стьюдента и Пирсона. Пользователям R нет нужды использовать таблицы, поскольку вычисления квантилей осуществляются функциями семейства **qname()**, где **name** — имя конкретного распределения в R, например для нормального распределения это функция **qnorm()**.

Если плотность вероятности оценок параметра неизвестна, но можно полагать, что она симметрична и унимодальна, то доверительные границы можно определить непосредственно на основе несмещенной точечной оценки $\hat{\theta}$ и ее среднеквадратичного отклонения $\sigma_{\hat{\theta}}$ по зависимостям

$$\hat{\theta}_H = \hat{\theta} - u_{y_H} \sigma_{\hat{\theta}} , \quad \hat{\theta}_B = \hat{\theta} + u_{y_B} \sigma_{\hat{\theta}} ,$$

где u_{y_H} и u_{y_B} это квантили нормированного центрированного распределения оценок, соответствующие вероятностям y_H и y_B . Зависимость обычно используется для интервального оценивания математического ожидания нормальной СВ с неизвестной дисперсией. При выборке ограниченного объема квантили определяются по распределению Стьюдента, а при большем объеме (как правило $> 50-100$) по нормальному распределению.

Квантиль — значение, которое СВ не превышает с заданной вероятностью, α -квантиль (т. е. квантиль порядка $\alpha \in (0; 1)$) это число x_α , такое что

$P\{X \leq x_\alpha\} \leq \alpha$, в случае же непрерывного распределения $F(x_\alpha) = \alpha$ (см. Рис.7).

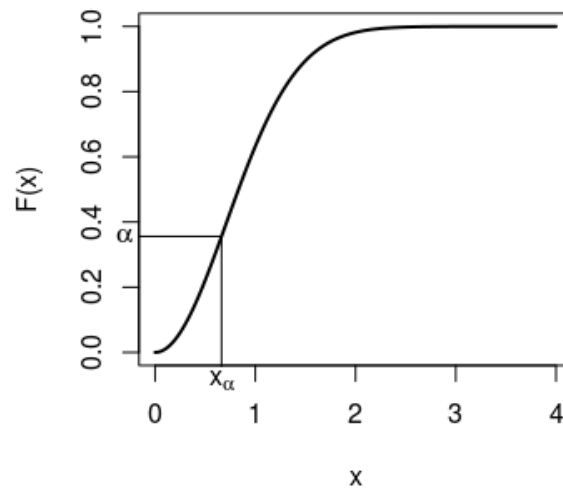


Рис. 7. Определение α -квантили.

Оценивание первых моментов распределения

Точечная оценка математического ожидания для выборки $x_i, i=1,2,\dots,n$ определяется как

$$\hat{M}_i = \frac{1}{n} \sum_{i=1}^n x_i .$$

Несмещенная оценка дисперсии:

$$\hat{\sigma}_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{M}_x)^2 .$$

Центральная предельная теорема утверждает, что распределение оценки МО имеет асимптотически нормальное распределение при $n \rightarrow \infty$ с точечной оценкой МО в качестве МО и дисперсией

$$\sigma_{\hat{M}_x}^2 = \frac{\sigma_x^2}{n} .$$

Интервальная оценка: $M \in [\hat{M}_x - u_p \sigma_{\hat{M}_x}; \hat{M}_x + u_p \sigma_{\hat{M}_x}]$, где u_p - квантиль распределения Стьюдента.

Проверка статистических гипотез

Статистическая гипотеза — предположение о типе неизвестного распределения и его параметрах. Проверяется на основе выборки реализаций СВ. Обычно используются следующие обозначения:

H_0 — нулевая гипотеза;

H_1 — конкурирующая (альтернативная) гипотеза, которая противоречит H_0 .

Гипотеза может быть *простой* (одна точечная оценка параметра) или *сложной* (множество значений, интервал, бесконечно много). Выдвинутая гипотеза может быть верной или ложной, проверка производится статистическими методами, оперирующими случайными величинами и поэтому не обеспечивают полной достоверности принимаемых решений. Могут быть следующие варианты, при принятии решений:

- 1) выдвигается верная гипотеза, и она принимается;
- 2) выдвигается верная гипотеза, но она отклоняется;
- 3) выдвигается ложная гипотеза, и она отклоняется;
- 4) выдвигается ложная гипотеза, но она принимается.

Ошибкой первого рода называется принятие альтернативной гипотезы H_1 , когда верна нулевая гипотеза H_0 . *Ошибкой второго рода* называется принятие гипотезы H_0 , когда верна альтернативная гипотеза H_1 . Возможные ошибки характеризуются их вероятностями. Вероятность совершить ошибку первого рода называется *уровнем значимости* и обозначается как α . Вероятность совершить ошибку второго рода принято обозначать как β . Вероятность отклонить нулевую гипотезу, когда верна альтернативная, называется *мощностью критерия*. Проверка гипотез осуществляется с помощью показателя Π , зависящего от элементов выборки, а различение гипотез производится благодаря тому, что априорные распределения показателя для H_0 и H_1 не совпадают. Правило принятия решения о верности той или иной гипотезы называют *статистическим критерием*, *статистическим тестом* или *критерием проверки статистической гипотезы*. На основе критерия определяется *критическая область*, попадание куда показателя Π означает то, что нулевая гипотеза отвергается. Критическая область может быть *двухсторонней* или *односторонней*. Напротив, *область принятия гипотезы* это совокупность значений показателя Π , при которых нулевая гипотеза принимается.

Гипотеза о параметре распределения первого типа — предположение, что неизвестный параметр θ равен θ_0 (простая нулевая гипотеза) или $\theta \in [\theta_1; \theta_2]$ (сложная нулевая гипотеза). Если получена оценка $\hat{\theta}$ с плотностью вероятности $f(\theta)$, то для проверки гипотезы задаются уровнем значимости α (вероятность ошибки первого рода), по выборке делается оценка доверительного интервала $[\hat{\theta}_H; \hat{\theta}_B]$ так что

$$P\{\theta \in [\hat{\theta}_H; \hat{\theta}_B]\} = 1 - \alpha = 1 - \alpha_1 - \alpha_2,$$

где $\alpha_1 = 1 - \gamma_H$, $\alpha_2 = 1 - \gamma_B$. Для проверки простой гипотезы $H_0: \theta = \theta_0$ используется решающее правило:

гипотеза принимается если $\theta_0 \in [\hat{\theta}_H; \hat{\theta}_B]$;

гипотеза отклоняется если $\theta_0 \notin [\hat{\theta}_H; \hat{\theta}_B]$.

Для проверки сложной нулевой гипотезы $H_0: \theta \in [\theta_1; \theta_2]$ применяется решающее правило:

гипотеза принимается если $[\hat{\theta}_H; \hat{\theta}_B] \in [\theta_1; \theta_2]$;

гипотеза отклоняется если $\hat{\theta}_H > \theta_2$ или $\hat{\theta}_B < \theta_1$.

В случае же частичного перекрытия отрезков, принять решение нельзя и необходимо уточнить гипотезу или увеличить объем выборки. Один из пределов множества $[\theta_1; \theta_2]$ может быть бесконечным или иметь естественные пределы, в этом случае оценивается односторонний доверительный интервал. В таком случае проверка гипотезы осуществляется по правилу:

гипотеза принимается, если $\hat{\theta}_B \leq \theta_2$ или $\hat{\theta}_H \geq \theta_1$;

гипотеза отклоняется в противном случае.

Гипотеза о параметрах распределения второго типа — предположение о том, что неизвестные параметры θ_1 и θ_2 в двух выборках совпадают, т. е.

$H_0: \theta_1 = \theta_2$. Для проверки такой гипотезы на основании выборок делаются оценки $\hat{\theta}_1$ и $\hat{\theta}_2$ и выбирается показатель $\Pi = F(\hat{\theta}_1, \hat{\theta}_2)$ (часто это просто $\hat{\theta}_1 - \hat{\theta}_2$ или $\hat{\theta}_1/\hat{\theta}_2$) и определяется его распределение как СВ. Для показателя вычисляется доверительный интервал при заданном уровне значимости. Решающее правило такое же как и для простой гипотезы первого типа, но для показателя Π .

Если имеется возможность построить распределения оценок $\hat{\theta}_1$ и $\hat{\theta}_2$, то можно непосредственно использовать вероятность неравенств $\theta_1 < \theta_2$ или $\theta_1 > \theta_2$. Для примера рассмотрим вариант плотностей вероятностей оценок, построенных по выборкам и показанный на Рис. 8.

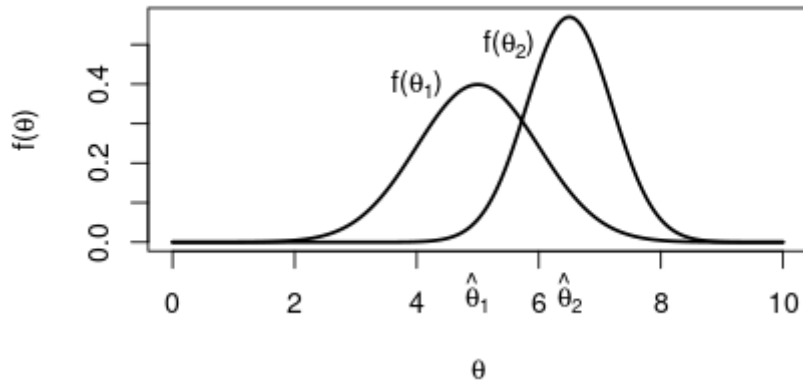


Рис.8. Пример плотности вероятности оценок

На основе плотностей вероятностей на Рис. 8 можно оценить вероятность того, что параметр θ_2 во второй выборке будет больше или меньше параметра θ_1 в первой выборке:

$$P\{\theta_2 > \theta_1\} = P_1 = \int_{-\infty}^{\infty} \int_{-\infty}^{\theta_1} f(\theta_1)f(\theta_2)d\theta_1d\theta_2,$$

$$P\{\theta_2 < \theta_1\} = P_2 = 1 - P_1,$$

$$P = \max[P_1; P_2].$$

В данном случае при рассматриваемом уровне значимости α критическая область для показателя (вероятности) P при проверке нулевой гипотезы $H_0: \theta_1 = \theta_2$ равна $P \in [1 - \alpha; 1]$. Поэтому решающее правило применяется следующее:

гипотеза принимается, если $P \leq 1 - \alpha$;

гипотеза отклоняется, если $P > 1 - \alpha$.

Гипотезы о распределении СВ, таким образом, делятся на два класса: гипотезы о совпадении распределений в двух выборках и гипотезы о типе распределения СВ. Для проверки гипотез используются специально подобранные показатели, отражающие степень близости распределений.

Методы оценивания параметров распределения СВ и проверки статистических гипотез

Для точечного оценивания параметров распределения СВ на основе выборок используются такие классические методы математической статистики как метод максимального правдоподобия (ММП), метод моментов (ММ) и метод наименьших квадратов (МНК).

Метод максимального правдоподобия

Для оценивания параметров распределения случайной величины по ММП необходимо знать тип распределения с точностью до неизвестного параметра $\theta = (\theta_1, \theta_2, \dots, \theta_J)$. Если СВ дискретная, $x_q, q=1, 2, \dots, Q$, то ее распределение описывается совокупностью вероятностей

$$P\{x_q | \theta\} = p_q, \quad q=1, 2, \dots, Q.$$

Если СВ непрерывная, то ее распределение можно представить в виде плотности вероятности или функции распределения:

$$f(x) = f(x, \theta), \quad F(x) = F(x, \theta).$$

Функция правдоподобия это априорная вероятность получения совокупности возможных значений СВ в выборке $\{x_i\}$, рассматриваемая как функция параметра θ . Для дискретной величины функция правдоподобия это

$$L(\theta) = P\{\{x_i\} | \theta\} = \prod_{i=1}^n P\{x_i | \theta\},$$

для непрерывной —

$$L(\theta) = f(\{x_i\}, \theta) = \prod_{i=1}^n f(x_i, \theta).$$

Сущность ММП в выборе параметра $\theta = (\theta_1, \theta_2, \dots, \theta_J)$, доставляющего максимум функции правдоподобия

$$L(\theta)|_{\theta=\hat{\theta}} = \max$$

или

$$\ln L(\theta)|_{\theta=\hat{\theta}} = \max.$$

Если ФП имеет аналитический вид, то оценка параметров может быть сделана на основе решения системы уравнений (необходимых условий экстремума)

$$\left. \frac{\partial L(\theta)}{\partial \theta_j} \right|_{\hat{\theta}_j} = 0, \quad j=1, 2, \dots, J$$

или

$$\left. \frac{\partial \ln L(\boldsymbol{\theta})}{\partial \theta_j} \right|_{\hat{\theta}_j} = 0, \quad j=1,2,\dots,J .$$

Точность получаемых оценок характеризуется *ковариационной матрицей точечных оценок*

$$\mathbf{K}_{\hat{\boldsymbol{\theta}}} = \begin{pmatrix} \sigma_{\hat{\theta}_1}^2 & K_{\hat{\theta}_1 \hat{\theta}_2} & \dots & K_{\hat{\theta}_1 \hat{\theta}_j} \\ K_{\hat{\theta}_2 \hat{\theta}_1} & \sigma_{\hat{\theta}_2}^2 & \dots & K_{\hat{\theta}_2 \hat{\theta}_j} \\ \dots & \dots & \dots & \dots \\ K_{\hat{\theta}_j \hat{\theta}_1} & K_{\hat{\theta}_j \hat{\theta}_2} & \dots & \sigma_{\hat{\theta}_j}^2 \end{pmatrix},$$

где $\sigma_{\hat{\theta}_j}^2$ - дисперсия j -ой компоненты вектора оценки параметра $\hat{\theta}_j$, а $K_{\hat{\theta}_j \hat{\theta}_k}$ - момент ковариации компонент вектора оценки $\hat{\theta}_j$ и $\hat{\theta}_k$. Аналитический вид матрицы $\mathbf{K}_{\hat{\boldsymbol{\theta}}}$ можно представить как обращение *информационной матрицы Фишера* при полученных точечных оценках

$$\mathbf{K}_{\hat{\boldsymbol{\theta}}} = - \left(\left. \frac{\partial^2 \ln L(\boldsymbol{\theta})}{\partial \theta_j \partial \theta_k} \right|_{\hat{\boldsymbol{\theta}}} \right)^{-1} .$$

Если θ скаляр, а не вектор, то дисперсия оценки определяется по более простой формуле

$$\sigma_{\hat{\theta}}^2 = - \left(\left. \frac{d^2 \ln L(\theta)}{d \theta^2} \right|_{\hat{\theta}} \right)^{-1} .$$

Доказано, что оценки ММП в общем случае являются асимптотически несмещенными и асимптотически эффективными (при объеме выборки $n \rightarrow \infty$).

Метод моментов

Рассмотрим СВ X с плотностью вероятности

$$f(x) = f(x, \boldsymbol{\theta}), \quad \boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_J) .$$

Для оценивания $\boldsymbol{\theta}$ надо получить J первых моментов. Для примера без потери общности положим $J = 2$, т. е.

$$f(x) = f(x, \theta_1, \theta_2) .$$

Два первых момента это математическое ожидание и дисперсия,

$$M_X = \int_{\Omega_x} x f(x, \theta_1, \theta_2) dx = F_1(\theta_1, \theta_2) ,$$

$$D_X = \int_{\Omega_x} (x - M_X)^2 f(x, \theta_1, \theta_2) dx = F_2(\theta_1, \theta_2) .$$

Если θ_1 и θ_2 неизвестны, но получена выборка $x_i, i=1,2,\dots,n$, то по ней можно сделать оценки \hat{M}_X и \hat{D}_X , дисперсии которых

$$\sigma_{\hat{M}_X}^2 = \frac{\hat{D}_X}{n}, \quad \sigma_{\hat{D}_X}^2 = \frac{2\hat{D}_X^2}{n} .$$

Точечные оценки параметров распределения θ_1 и θ_2 вычисляются путем приравнивания выражений для моментов распределения их точечным оценкам. Т.е. для случая $J = 2$ нужно решить следующую систему уравнений:

$$\begin{cases} F_1(\hat{\theta}_1, \hat{\theta}_2) = \hat{M}_X \\ F_2(\hat{\theta}_1, \hat{\theta}_2) = \hat{D}_X \end{cases} .$$

Решение этой системы выполняется различными методами в зависимости от сложности функций F_1 и F_2 . Иногда в простых случаях можно получить явные выражения вида

$$\theta_j = \varphi_j(\hat{M}_X, \hat{D}_X), \quad j=1,2 .$$

Дисперсии оценок с использованием линеаризации определяются как

$$\sigma_{\hat{\theta}_j}^2 = \left[\frac{\partial \varphi_j(\hat{M}_X, \hat{D}_X)}{\partial \hat{M}_X} \right]^2 \sigma_{\hat{M}_X}^2 + \left[\frac{\partial \varphi_j(\hat{M}_X, \hat{D}_X)}{\partial \hat{D}_X} \right]^2 \sigma_{\hat{D}_X}^2, \quad j=1,2 .$$

Пример: Оценивание параметров по ММ равномерного распределения

$$f(x) = \begin{cases} \frac{1}{B-A}, & x \in [A; B] \\ 0, & x \notin [A; B] \end{cases} .$$

Для получения оценок параметров A и B нужно решить систему

$$\begin{cases} \frac{1}{2}(\hat{A} + \hat{B}) = \hat{M}_x \\ \frac{(\hat{B} - \hat{A})^2}{12} = \hat{D}_x \end{cases},$$

решение которой

$$\begin{cases} \hat{A} = \hat{M}_x - \sqrt{3\hat{D}_x} \\ \hat{B} = \hat{M}_x + \sqrt{3\hat{D}_x} \end{cases}.$$

Выражения для дисперсий оценок параметров A и B имеют вид

$$\sigma_{\hat{A}}^2 = \sigma_{\hat{B}}^2 = \sigma_{\hat{M}_x}^2 + \frac{3\sigma_{\hat{D}_x}^2}{4\hat{D}_x}.$$

Нормальное распределение имеет два параметра: математическое ожидание и дисперсию. Поэтому нет нужды решать систему уравнение по ММ, поскольку оценки ММ совпадают с точечными оценками по выборке,

$$\hat{M}_x = \frac{1}{n} \sum_{i=1}^n x_i, \quad \sigma_{\hat{M}_x}^2 = \frac{\hat{D}_x}{n};$$

$$\hat{D}_x = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{M}_x)^2, \quad \sigma_{\hat{D}_x}^2 = \frac{2\hat{D}_x^2}{n}.$$

Метод наименьших квадратов

Метод наименьших квадратов является распространенным методом оценивания параметров любых параметрических моделей, не обязательно вероятностных распределений. Его можно использовать для оценки первых двух моментов по стабильной или динамичной выборке. МНК можно использовать для решения уравнений, связанных с ММП и ММ. МНК так же используется для оценивания динамики неслучайных параметров и функций по результатам измерений, искаженных случайными погрешностями, а также для поиска зависимостей исследуемых показателей от различных измеряемых параметров и факторов.

В качестве примера рассмотрим применение МНК для подбора параметров модели описывающей динамику математического ожидания:

$$M(t) = F(t, \theta),$$

где $F()$ это известная функция или алгоритм, а $\theta = (\theta_1, \theta_2, \dots, \theta_J)$ - вектор неизвестных параметров. Наблюдения $x_i, i=1, 2, \dots, n$ в моменты времени t_1, t_2, \dots, t_n интерпретируются как

$$x_i = F(t_i, \theta) + \delta_i ,$$

где δ_i - случайные флуктуации, как правило с нулевым средним нормальным распределением. Собственно МНК в общем виде это решение оптимизационной задачи вида

$$\min_{\theta} \sum_{i=1}^n (x_i - F(t_i, \theta))^2 ,$$

закключающейся в поиске параметров θ , которые доставляют минимум сумме квадратов разностей между измеряемой величиной и функцией ее моделирующей. В общем случае задача может оказаться сложной и не иметь аналитического решения. Задача может оказаться плохо обусловленной и требовать хорошего априорного приближения вектора параметров θ_0 .

Для упрощения решения задачи, модель может представляться в *линеаризованном* виде:

$$M(t) \approx F(t, \theta_0) + \sum_{j=1}^J \left. \frac{\partial F(t, \theta)}{\partial \theta_j} \right|_{\theta_0} (\theta_j - \theta_{j0}) = F(t, \theta_0) + \sum_{j=1}^J \varphi_j(t) \Delta \theta_j .$$

Модель наблюдений уже записывается в линеаризованном виде следующим образом:

$$\sum_{j=1}^J \varphi_j(t_i) \Delta \theta_j + \delta_i = \Delta x_i , \quad i=1, 2, \dots, n ,$$

где $\Delta x_i = x_i - F(t_i, \theta_0)$. Собственно в данной записи линеаризованная модель это система линейных уравнений.

Линеаризованную модель можно записать в вектор-матричной форме как

$$\Phi \cdot \Delta \Theta + \Delta = \Delta X ,$$

где $\Delta \theta = \{\theta_j\}$ - вектор отклонений неизвестных параметров от начального приближения, $\Phi = \{\varphi_j(t_i)\}$ - матрица преобразования, Δ - вектор случайных флуктуаций, ΔX - вектор отклонений данных наблюдений от начальных значений модели.

Если известна ковариационная матрица K_{Δ} случайного вектора Δ , то обобщенный метод наименьших квадратов предполагает минимизацию следующей целевой функции:

$$S_{об} = V^T K_{\Delta}^{-1} V = \min_{\theta},$$

где $V = (x_1 - F(\theta, t_1), x_2 - F(\theta, t_2), \dots, x_n - F(\theta, t_n))^T$ - вектор отклонений.

Если флуктуации не коррелированы и имеют разные дисперсии, то используется взвешенный метод наименьших квадратов, предполагающий минимизацию следующей целевой функции:

$$S_{взв} = \sum_{i=1}^n \sigma_i^2 [x_i - F(\theta, t_i)]^2 = \min.$$

И наконец, если флуктуации не коррелированы и имеют постоянную дисперсию, т. е. $\sigma_i = const = \sigma, \forall i$, то применяется классический МНК с целевой функцией вида

$$S_{кл} = \sum_{i=1}^n [x_i - F(\theta, t_i)]^2 = \min.$$

Для линейной или линеаризованной модели наблюдений аналогичные целевые функции легко получить подстановкой вместо $F(\theta, t)$ её линеаризованной версии. Преимущество линейного варианта заключается в том, что оценки отклонений параметра $\Delta \hat{\theta}$ линейно зависят от ΔX , т. е. $\Delta \hat{\theta} = B \Delta X$, где B это некоторая матрица преобразования и в этом случае легко получить ковариационную матрицу оценки векторного параметра:

$$K_{\Delta \hat{\theta}} = B K_{\Delta} B^T.$$

Для линеаризованной модели наблюдений для вычисления точечных оценок и ковариационной матрицы оценок вектора параметров готовые выражения имеют следующий вид:

для обобщенного МНК

$$\Delta \hat{\theta} = (\Phi^T K_{\Delta}^{-1} \Phi)^{-1} \Phi^T K_{\Delta}^{-1} \Delta X, \quad K_{\Delta \hat{\theta}} = (\Phi^T K_{\Delta}^{-1} \Phi)^{-1};$$

для взвешенного МНК

$$\Delta \hat{\theta} = (\Phi_n^T \Phi_n)^{-1} \Phi_n^T \Delta X_n, \quad K_{\Delta \hat{\theta}} = (\Phi_n^T \Phi_n)^{-1} ;$$

для классического МНК

$$\Delta \hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T \Delta X, \quad K_{\Delta \hat{\theta}} = \sigma^2 (\Phi^T \Phi)^{-1} .$$

Здесь в формулах для взвешенного МНК использованы взвешенные значения вектора и матрицы, т. е. $\Delta X = [\Delta x_i / \sigma_i]$, $\Phi_n = [\varphi_j(t_i) / \sigma_i]$. Если в случае классического МНК дисперсия флуктуаций σ^2 неизвестна, то *несмещенную точечную оценку дисперсии флуктуаций* можно определить по остаточным невязкам как

$$\hat{\sigma}^2 = \frac{1}{n-J} (\Delta X - \Phi \Delta \hat{\theta})^T (\Delta X - \Phi \Delta \hat{\theta}) .$$

На основе оценок вектора отклонений вектора параметров от начального приближения определяется оценка самого вектора и его ковариационной матрицы:

$$\hat{\theta} = \theta_0 + \Delta \hat{\theta}; \quad K_{\hat{\theta}} = K_{\Delta \hat{\theta}} .$$

Регрессионный анализ

При анализе данных наблюдений могут решаться задачи исследования зависимости параметров распределения СВ от различных факторов. Для таких задач используются методы *регрессионного анализа*.

Рассмотрим зависимость вида

$$M_X(\mathbf{g}) = a_0 + a_1 g_1 + a_2 g_2 + \dots + a_J g_J = a_0 + \sum_{j=1}^J a_j g_j = a_0 + \mathbf{a}^T \mathbf{g} .$$

Для простоты записи можно полагать, что всегда $g_0 = 1$ и тогда запись упростится, т. е.

$$M_X(\mathbf{g}) = \mathbf{a}^T \mathbf{g} ,$$

где векторы $\mathbf{a} = (a_0, a_1, a_2, \dots, a_J)^T$, $\mathbf{g} = (1, g_1, g_2, \dots, g_J)^T$, размерность которых $J+1$.

Если от какого-то фактора имеется нелинейная зависимость, например

$$M_x(z) = b_1 z + b_2 z^2 ,$$

то z и z_2 представляются как отдельные факторы, т. е. $z = g_k$, $z^2 = g_{k+1}$ и можно снова пользоваться аппаратом линейной регрессии.

Для оценивания коэффициентов регрессии a_0, a_1, \dots, a_J проводятся наблюдения с выборкой $x_i, i=1, 2, \dots, n$ при различных значениях факторов g_1, g_2, \dots, g_J (т. е. вместе с выборкой нужно получить набор $\mathbf{g}_i, i=1, 2, \dots, n$) и используется следующая модель наблюдений:

$$M_x(\mathbf{g}) + \delta_i = x_i, \quad i=1, \dots, n$$

Задача регрессионного анализа сводится к МНК для решения задачи нахождения минимума (в пространстве коэффициентов регрессии)

$$\min_{\mathbf{a}} \sum_{i=1}^n \left(x_i - \sum_{j=0}^J a_j g_{ji} \right)^2 ,$$

здесь надо помнить о соглашении $g_{0i} = 1, \forall i$. Задача может быть переформулирована в векторно-матричном виде как

$$\min_{\mathbf{a}} (\mathbf{x} - \mathbf{G}\mathbf{a})^T (\mathbf{x} - \mathbf{G}\mathbf{a}) .$$

В соответствии с классическим МНК можно получить зависимости для точечных оценок вектора регрессионных коэффициентов и их ковариационной матрицы:

$$\hat{\mathbf{a}} = (\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \mathbf{x} ;$$

$$\hat{\sigma}^2 = \frac{1}{n-J} (\mathbf{x} - \mathbf{G}\hat{\mathbf{a}})^T (\mathbf{x} - \mathbf{G}\hat{\mathbf{a}}) ;$$

$$\mathbf{K}_{\hat{\mathbf{a}}} = \hat{\sigma}^2 (\mathbf{G}^T \mathbf{G})^{-1} .$$

Все эти вычисления возможны только если все факторы линейно независимы, если какие-то факторы линейно зависимы, то необходимо уменьшить число этих факторов, оставив только линейно независимые или функции от них.

Упражнение. Выполнить работу предложенную по ссылке
<http://environmentalcomputing.net/linear-regression/>

Методы интервального оценивания параметров распределения СВ

Оперативный метод

Если для неизвестного параметра θ получена точечная оценка $\hat{\theta}$ и ее среднеквадратичное отклонение $\sigma_{\hat{\theta}}$, то оценка границ доверительного интервала $[\theta_H; \theta_B]$ при заданной доверительной вероятности γ

$$\hat{\theta}_H = \hat{\theta} - u_{\gamma_H} \sigma_{\hat{\theta}},$$

$$\hat{\theta}_B = \hat{\theta} + u_{\gamma_B} \sigma_{\hat{\theta}},$$

где $\gamma_H = P\{\theta \geq \theta_H\}$, $\gamma_B = P\{\theta \leq \theta_B\}$, $\gamma_H + \gamma_B = 1 + \gamma$, $u_{\gamma_H}, u_{\gamma_B}$ - квантили распределения точечной оценки, соответствующие вероятностям γ_H и γ_B соответственно. Если при решении практических задач предполагают нормальное распределение оценки параметра $\hat{\theta}$, то квантили распределения определяются по формуле Лапласа. В R это `qnorm(γ, 0, 1)`, т. е.

$$u_{\gamma} = \text{qnorm}(\gamma, 0, 1).$$

Если распределение $\hat{\theta}$ отличается от нормального, но известно, что оно симметрично и унимодально, то квантили больше

$$u_{\gamma} = \frac{\sqrt{2}}{3\sqrt{1-\gamma}}.$$

Если о распределении вообще ничего не известно, то квантили еще больше

$$u_{\gamma} = \frac{1}{\sqrt{2(1-\gamma)}}.$$

Упражнение в R: Пусть $\hat{M}_x = 3.0$, $\hat{\sigma}_{\hat{M}_x} = 0.45$, $\gamma = 0.9$. Определить $[\theta_H; \theta_B]$ для всех вышеописанных случаев.

Интервальное оценивание МО СВ с нормальным распределением и неизвестной дисперсией

В этом случае используется показатель (статистика)

$$t = \frac{\hat{M}_X - M_X}{\sigma_{\hat{M}_X}}$$

имеющий распределение Стьюдента с числом степеней свободы $k=n-1$, где n это объем выборки. (см. презентацию «Распределения дискретных и непрерывных случайных величин и знакомство с ними в R»). Оценка границ доверительного интервала выполняется по формулам

$$\hat{M}_{XH} = \hat{M} - t(\gamma, k) \sigma_{\hat{M}_X} ,$$

$$\hat{M}_{XB} = \hat{M} + t(\gamma, k) \sigma_{\hat{M}_X} ,$$

где $t(\gamma, k)$ - квантиль распределения Стьюдента, которая в R вычисляется с помощью функции $qt(\gamma, k, ncp=0)$.

Связь интервальных и точечных оценок

Распределение Стьюдента получено для случайной величины t при фиксированном значении M_X , для случайной величины X , имеющей нормальное распределение. Если распределение Стьюдента верно отражает распределение t , то можно построить распределение \hat{M}_X . Обозначим неизвестную величину M_X как m , тогда

$$t = \frac{\hat{M}_X - m}{\sigma_{\hat{M}_X}} .$$

Тогда плотность распределения m может быть выражена через распределение Стьюдента как

$$g(m) = f\left(\frac{\hat{M}_X - m}{\sigma_{\hat{M}_X}}\right) \cdot \left|\frac{dt}{dm}\right| = \frac{\Gamma((k+1)/2)}{\sqrt{k\pi}\Gamma(k/2)} \cdot \left[1 + k^{-1} \left(\frac{\hat{M}_X - m}{\sigma_{\hat{M}_X}}\right)^2\right]^{-\frac{k+1}{2}} \cdot \left(\frac{1}{\sigma_{\hat{M}_X}}\right)$$

Интегрирование этой плотности приводит к тем же оценкам границ доверительного интервала $[\hat{M}_{XH}; \hat{M}_{XB}]$, а также можно воспользоваться равенствами

$$\hat{M}_X = \int_{-\infty}^{\infty} m g(m) dm ,$$

$$\sigma_{\hat{M}_X}^2 = \int_{-\infty}^{\infty} (m - \hat{M}_X)^2 g(m) dm .$$

Оценка математического ожидания совпадает с точечной оценкой по выборке, т. е. с $\frac{1}{n} \sum_{i=1}^n x_i$, однако дисперсия отличается от дисперсии определяемой по ММП.

Определение доверительных границ для дисперсии

Рассмотрим случайную величину X с неизвестными M_X и D_X . Пусть имеется выборка $\{x_i\}$ объема n . Используется в данном случае показатель (статистика) следующего вида:

$$\chi^2 = \frac{n\hat{D}_X}{D_X} .$$

При нормальном распределении случайной величины X показатель χ^2 имеет χ^2 -распределение с числом степеней свободы $k = n - 1$. Для удобства записи введем обозначение $y = \chi^2 \in [0; +\infty]$, тогда функция плотности вероятности имеет вид

$$f(y) = [\Gamma(k/2) \cdot 2^{k/2}]^{-1} \cdot y^{(k/2-1)} \cdot e^{-y/2} ,$$

в R это функция `dchisq()`.

Для оценки границ доверительного интервала задаются вероятностями γ_H и γ_B так чтобы

$$P\{D \leq D_B\} = 1 - \gamma_H ,$$

$$P\{D \geq D_H\} = \gamma_B .$$

Затем определяются квантили $\chi_1^2(\gamma_H, k)$ и $\chi_2^2(\gamma_B, k)$ и определяются доверительные границы по формулам

$$\hat{D}_{XH} = \frac{n\hat{D}_X}{\chi_2^2(\gamma_B, k)} ,$$

$$\hat{D}_{XB} = \frac{n\hat{D}_X}{\chi_1^2(\gamma_H, k)} .$$

Пример: $n = 10$, $\hat{D}_X = 2.03$, $\gamma_H = 0.05$, $\gamma_B = 0.95$. Можно получить для этого случая

$$\hat{D}_{XH} \approx 1.20, \quad \hat{D}_{XB} \approx 6.11 .$$

Упражнение в R: получить ответы примера самостоятельно используя функцию `qchisq()`.

Связь интервальных и точечных оценок

Построим на основе χ^2 -распределения для статистики χ^2 распределение для D_X .

$$y = \frac{n\hat{D}_X}{d} ,$$

$$g(d) = f\left(\frac{n\hat{D}_X}{d}\right) \left|\frac{\partial y}{\partial d}\right| = [\Gamma(k/2) 2^{k/2}]^{-1} \left(\frac{n\hat{D}_X}{d}\right)^{k/2-1} e^{-\frac{n\hat{D}_X}{2d}} \cdot \left(\frac{n\hat{D}_X}{d^2}\right) .$$

В правой части все выражение, за исключением последней скобки вычисляется в R функцией `dchisq(y, k, ...)`. На основе этой плотности можно вычислить несмещенную точечную оценку дисперсии и ее дисперсию,

$$\hat{D} = \int_0^{\infty} d \cdot g(d) \cdot dd; \quad \sigma_D^2 = \int_0^{\infty} (d - \hat{D})^2 g(d) dd .$$

Эти оценки не совпадают с оценками ММП.

Методы проверки статистических гипотез

Проверка гипотезы о математическом ожидании случайной величины

Пусть наша нулевая гипотеза заключается в том, что математическое ожидание случайной величины X совпадает с некоторым значением M_0 , т. е.

$$H_0: M_X = M_0 ;$$

$$H_1: M_X \neq M_0 .$$

Проверка гипотезы осуществляется следующим образом. Задается уровень значимости α , вычисляются выборочные оценки

$$\hat{M}_X = \frac{1}{n} \sum_{i=1}^n x_i; \quad \sigma_{\hat{M}_X}^2 = \frac{\sigma_X^2}{n} .$$

Если дисперсия σ_X^2 неизвестна, то делается и ее выборочная оценками

$$\hat{\sigma}_X^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{M}_X)^2 .$$

Для проверки гипотезы используется следующий показатель (статистика):

$$z = \frac{\hat{M}_X - M_0}{\sigma_{\hat{M}_X}} .$$

В зависимости от того известна ли σ_X^2 , используются два способа.

1-ый способ, σ_X^2 известна. В этом случае показатель описывается стандартным нормальным распределением с нулевым средним и единичной дисперсией, т. е. $z \in N[0,1]$. Определяется критическое значение показателя $z_{кр}$ на основе функции распределения стандартного нормального распределения (с нулевым средним и единичной дисперсией) с использованием функции Лапласа из соотношения

$$\Phi(z_{кр}) = 1 - \alpha .$$

В учебной литературе по математической статистике имеется разночтение, что именно считать функцией Лапласа, это либо $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{z^2}{2}} dz$, либо

$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{z^2}{2}} dz$ (выражения отличаются нижним пределом интегрирования), поэтому уравнение для определения $z_{кр}$ может быть записано как

$$0.5 + \Phi(z_{кр}) = 1 - \alpha ,$$

но в данном курсе мы будем считать, что функция $\Phi(\cdot)$ это функция распределения стандартного нормального распределения. Таким образом $z_{кр}$ можно определить функцию вычисления квантилей в R

$$z_{кр} = qnorm(1 - \alpha, 0, 1) .$$

По выборке вычисляется показатель \hat{z} и сравнивается с $z_{кр}$. Для проверки гипотезы о математическом ожидании случайной величины используется следующее решающее правило:

гипотеза $H_0: M_X = M_0$ принимается, если $|\hat{z}| \leq z_{кр}$,

гипотеза H_0 отклоняется, если $|\hat{z}| > z_{кр}$.

2-ой способ, σ_X^2 неизвестна. Вместо дисперсии используется ее выборочная оценка $\hat{\sigma}_X^2$ и вычисляются уже квантили распределения Стьюдента с $k = n - 1$ числом степеней свободы, т. е. в R это

$$z_{кр} = qt(1 - \alpha, n - 1).$$

Правило проверки гипотезы то же, что и для 1-го способа.

Следует отметить, что назначение слишком низкого уровня значимости может привести к отбрасыванию альтернативной гипотезы, когда она все таки верна. Поэтому назначение низкого уровня значимости следует делать тогда, когда есть высокая уверенность в верности нулевой гипотезы. В противном случае при малом уровне α с большой вероятностью $1 - \alpha$ может быть отклонена верная альтернативная гипотеза.

Проверка гипотез о дисперсии случайной величины

Пусть имеется выборка $x_i, i = 1, 2, \dots, n$ из генеральной совокупности с неизвестными математическим ожиданием M_X и дисперсией D_X . Проверяется гипотеза о том, что дисперсия данной СВ не превышает некоторого заданного значения D_0 , т. е. проверяются гипотезы

$$H_0: D_X \leq D_0,$$

$$H_1: D_X > D_0.$$

Если СВ имеет нормальное распределение, то проверка H_0 делается по критерию χ^2 К.Пирсона

$$\chi^2 = (n - 1) \frac{\hat{\sigma}_X^2}{D_0},$$

где $\hat{\sigma}_X^2 = \frac{1}{n - 1} \sum_{i=1}^n (x_i - \hat{M}_X)^2$, $\hat{M}_X = \frac{1}{n} \sum_{i=1}^n x_i$. Задается α - уровень значимости, при нормальном распределении СВ, показатель имеет χ^2 -

распределение с $k=n-1$ числом степеней свободы. Затем находится квантиль $\chi_{кр}^2$ такой что

$$P\{\chi^2 > \chi_{кр}^2\} = \alpha ,$$

в R это делается легко $\chi_{кр}^2 = qchisq(1-\alpha, k)$. Применяется решающее правило:

гипотеза $H_0: D_X \leq D_0$ принимается, если $\chi^2 \leq \chi_{кр}^2$;

гипотеза H_0 отклоняется если $\chi^2 > \chi_{кр}^2$.

Принятие нулевой гипотезы лишь означает, что данные наблюдений не противоречат ей. При назначении α следует учитывать, что вероятность отклонить альтернативную гипотезу, если она верна, равна $1-\alpha$.

При проверке гипотез полезно определить вероятность того, что при справедливой нулевой гипотезе значение показателя χ^2 может превышать полученную реализацию $\hat{\chi}^2$. Такая вероятность определяется на основе выражения

$$P\{\chi^2 > \hat{\chi}^2\} = 1 - \gamma(\hat{\chi}^2) ,$$

где $\gamma(\hat{\chi}^2)$ - вероятность, определяемая по функции χ^2 -распределения, в R это делается с помощью функции $pchisq(\hat{\chi}^2, k)$. Эта величина показывает вероятность того, что при справедливости нулевой гипотезы оценка дисперсии по выборке могла быть еще выше.

Проверка гипотезы о математических ожиданиях случайных величин по двум выборкам

Пусть имеется две выборки $\{x_i\}$ и $\{y_i\}$ объемами n_1 и n_2 соответственно. Рассмотрим гипотезу о равенстве M_X и M_Y для случая нормального распределения случайных величин в обоих выборках. Для проверки гипотезы используются точечные оценки

;

$$\hat{M}_Y = \frac{1}{n_2} \sum_{i=1}^{n_2} y_i ; \quad \sigma_{\hat{M}_Y}^2 = \frac{\sigma_Y^2}{n_2} .$$

Используется показатель $z = \frac{\hat{M}_X - \hat{M}_Y}{\sigma_{\Delta M}}$, где $\sigma_{\Delta M}$ - среднеквадратическое отклонение разности разности оценок математических ожиданий:

$$\sigma_{\Delta M} = \sqrt{\frac{\sigma_X^2}{n_1} + \frac{\sigma_Y^2}{n_2}} .$$

Способ проверки зависит от того, известны ли дисперсии СВ или нет.

1-ый способ, дисперсии σ_X^2 и σ_Y^2 известны. В этом случае при нормально распределенных случайных величинах, показатель z имеет стандартное нормальное распределение (с нулевым средним и единичной дисперсией), поэтому для проверки гипотезы задается уровень значимости α , определяется $z_{кр}$ (с использованием R)

$$z_{кр} = qnorm(1-\alpha, 0, 1)$$

и используется решающее правило

гипотеза $H_0: M_X = M_Y$ принимается, если $|\hat{z}| \leq z_{кр}$;

гипотеза H_0 отклоняется, если $|\hat{z}| > z_{кр}$.

2-ой способ, дисперсии σ_X^2 и σ_Y^2 известны. Вместо известных дисперсий в этом случае используются их оценки

$$\hat{\sigma}_X^2 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (x_i - \hat{M}_X)^2 ; \quad \hat{\sigma}_Y^2 = \frac{1}{n_2-1} \sum_{i=1}^{n_2} (y_i - \hat{M}_Y)^2 ,$$

а проверка гипотезы проводится по t -критерию Стьюдента

$$t = \frac{|\hat{M}_X - \hat{M}_Y|}{\hat{\sigma}_{\Delta M}} , \text{ где } \hat{\sigma}_{\Delta M} = \sqrt{\frac{\hat{\sigma}_X^2}{n_1} + \frac{\hat{\sigma}_Y^2}{n_2}} .$$

Определяется критическое значение критерия, в R используется функция qt():

$$t_{кр} = qt(1-\alpha, n_1+n_2-2) .$$

Решающее правило в этом случае:

гипотеза $H_0: M_X = M_Y$ принимается, если $|\hat{t}| \leq t_{кр}$;

гипотеза H_0 отклоняется, если $|\hat{t}| > t_{кр}$.

Проверка гипотез о дисперсиях СВ по двум выборкам

Итак, имеется две выборки $x_i, i=1,2,\dots,n_1$ и $y_i, i=1,2,\dots,n_2$ и нужно проверить гипотезу о равенстве дисперсий, т. е. $H_0: \sigma_X^2 = \sigma_Y^2$. Используется F -критерий Фишера-Снедекора

$$F = \frac{\hat{\sigma}_X^2}{\hat{\sigma}_Y^2},$$

который при нормальных СВ и справедливости нулевой гипотезы имеет F -распределение Фишера-Снедекора, которое зависит от двух чисел степеней свободы $k_1=n_1-1$ и $k_2=n_2-1$. Выбирается критическое значение F из условия

$$P\{F > F_{кр}\} = \alpha,$$

что в среде R сделать просто с помощью функции $qf()$

$$F_{кр} = qf(1-\alpha, k_1, k_2, ncp=0).$$

И наконец, применяется решающее правило:

гипотеза $H_0: \sigma_X^2 = \sigma_Y^2$ принимается, если $\hat{F} \leq F_{кр}$;

гипотеза отклоняется, если $\hat{F} > F_{кр}$.

Проверка гипотезы об однородности распределений СВ в двух выборках

Гипотезу можно сформулировать как $H_0: f(x) = f(y)$. Наиболее мощным для этого случая считается критерий Вилкоксона (Wilcoxon), который применяется к упорядоченным в виде вариационного ряда выборкам

$$\{x_i\} = x_1, x_2, \dots, x_{n_1}, \quad x_1 \leq x_2 \leq x_3 \leq \dots \leq x_{n_1},$$

$$\{y_i\} = y_1, y_2, \dots, y_{n_2}, \quad y_1 \leq y_2 \leq y_3 \leq \dots \leq y_{n_2}.$$

Выборки объединяются в общий вариационный ряд, например такой

$$\{x, y\} = x_1, x_2, y_1, x_3, y_2, \dots, y_{n_2}, x_{n_1},$$

элементы которого нумеруются сквозным образом, а номер (или ранг) перечисляет элементы общего вариационного ряда

$$r_q = 1, 2, 3, 4, 5, \dots, n-1, n ,$$

где $n = n_1 + n_2$. Вводятся показатели следующего вида

$$U_X = \sum_{\text{только для } \{x_i\}} r_{qx} ,$$

$$U_Y = \sum_{\text{только для } \{y_i\}} r_{qy} ,$$

в которых просуммированы только номера из общего вариационного ряда, соответствующие только одной выборке. Если нулевая гипотеза верна, то при $n_1 + n_2 \geq 20$ и $\min(n_1, n_2) \geq 5$ распределение критерия Вилкоксона близко к нормальному со следующими математическими ожиданиями и дисперсией:

$$M_{U_x} = \frac{1}{2} n_1 (n_1 + n_2 + 1) ,$$

$$M_{U_y} = \frac{1}{2} n_2 (n_1 + n_2 + 1) ,$$

$$\sigma_U^2 = \frac{1}{12} n_1 n_2 (n_1 + n_2 + 1) .$$

Для проверки нулевой гипотезы определяются границы критической области для показателей

$$U_{HX} = M_{U_x} - u_\gamma \sigma_U ; \quad U_{BX} = M_{U_x} + u_\gamma \sigma_U ;$$

$$U_{HY} = M_{U_y} - u_\gamma \sigma_U ; \quad U_{BY} = M_{U_y} + u_\gamma \sigma_U ,$$

где u_γ это квантиль стандартного нормального распределения для $\gamma = 1 - \frac{\alpha}{2}$ (в R $u_\gamma = qnorm(\gamma, 0, 1)$). Решающее правило формулируется следующим образом:

гипотеза $H_0: f(x) = f(y)$ принимается, если $U_X \in [U_{HX}; U_{BX}]$ и $U_Y \in [U_{HY}; U_{BY}]$;

иначе гипотеза H_0 отклоняется.

Проверка гипотез о типе распределения

При априорной неизвестности типа распределения СВ, он наиболее достоверно определяется на основе гистограммы распределения по χ^2 -критерию К.Пирсона или на основе эмпирической функции распределения по ω^2 -критерию Крамера-Мизеса-Смирнова.

χ^2 -критерий

Пусть имеется выборка $x_i, i=1,2,\dots,n$ и предполагается распределение СВ в виде $f(x, \theta)$, где $\theta = (\theta_1, \theta_2, \dots, \theta_J)$. Строится гистограмма из K отрезков длиной Δ , определяются числа n_k элементов выборки попавших в k -ый отрезок. Теоретическая вероятность попадания в k -ый отрезок с выбранной функцией плотности вероятности

$$p_{kT} = f(x_k, \theta) \cdot \Delta, \quad k=1,2,\dots,K,$$

где x_k - средняя точка k -ого отрезка гистограммы. Параметры θ оцениваются на основе статистических методов. С другой стороны, на основе выборки легко определяется экспериментальная вероятностью

$$p_{k\text{эксн}} = \frac{n_k}{n}, \quad k=1,2,\dots,K.$$

Используется критерий

$$\chi^2 = n \sum_{k=1}^K \frac{(p_{kT} - p_{k\text{эксн}})^2}{p_{kT}},$$

который при нормальном распределении СВ имеет точно χ^2 -распределение с $q = K - J - 1$ степенями свободы. Доказано, что в пределе, при $n \rightarrow \infty$ распределение показателя стремится к χ^2 -распределению и при других распределениях СВ. Проверка гипотезы осуществляется следующим образом. Выбирается уровень значимости α и определяется критический показатель из условия $P\{\chi^2 \geq \chi_{\text{кр}}^2\} = \alpha$, что в R можно сделать через вызов функции `qchisq()`, т. е.

$$\chi_{\text{кр}}^2 = \text{qchisq}(1 - \alpha, q).$$

Затем вычисляется реализация $\chi_{\text{эксн}}^2$ и сравнивается с критическим показателем, т. е. решающее правило имеет вид:

гипотеза $H_0: f(x) = f(x, \theta)$ принимается, если $\chi_{\text{эксн}}^2 < \chi_{\text{кр}}^2$;

гипотеза H_0 отклоняется, если $\chi_{\text{эксп}}^2 \geq \chi_{\text{кр}}^2$.

ω^2 -критерий Крамера-Мизеса-Смирнова

Данный критерий основан на непосредственном сравнении функции распределения $F(x, \theta)$ с эмпирической ФР, основанной на вариационном ряде. В качестве меры отклонения гипотетической ФР от эмпирической используется средний квадрат разности функций распределения по всем значениям аргумента:

$$\omega^2 = n \int_{-\infty}^{\infty} [F_{\text{э}}(x) - F(x)]^2 dF(x) .$$

После преобразования интеграла к сумме, со значениями в точках выборки, с учетом зависимости эмпирической ФР от упорядоченных элементов выборки, показатель принимает вид

$$\omega^2 = s = \frac{1}{12n} + \sum_{i=1}^n \left[F(x_i) - \frac{2i-1}{2n} \right]^2 .$$

Математическое ожидание и дисперсия этого показателя

$$M_{\omega^2} = \frac{1}{6} ; \quad \sigma_{\omega^2}^2 = \frac{4n-3}{180n} .$$

Известно, что при $n > 40$ распределение этого показателя близко к предельному, которое можно использовать для вычисления квантилей (хотя можно найти таблицы). Предельное распределение имеет вид

$$a_1(s) = \frac{1}{\sqrt{2s}} \sum_{j=0}^{\infty} \frac{\Gamma(j+1/2)\sqrt{4j+1}}{\Gamma(1/2)\Gamma(j+1)} \exp\left\{-\frac{(4j+1)^2}{16s}\right\} \left\{ I_{-\frac{1}{4}}\left(\frac{(4j+1)^2}{16s}\right) - I_{\frac{1}{4}}\left(\frac{(4j+1)^2}{16s}\right) \right\} ,$$

где $I_{-\frac{1}{4}}$ и $I_{\frac{1}{4}}$ - модифицированные функции Бесселя, которые можно представить в виде следующего выражения

$$I_{\nu}(z) = \sum_{k=0}^{\infty} \frac{\left(\frac{z}{2}\right)^{\nu+2k}}{\Gamma(k+1)\Gamma(k+\nu+1)} , \quad |z| < \infty, \quad |\arg z| < \pi .$$