

Введение

Статистический анализ данных (САД) – это научная дисциплина, основанная на методах теории вероятностей и математической статистики, целью которой является формирование научно обоснованных выводов и принятие решений относительно сложной системы (объекта или явления) на основе статистических данных о ней.

Введем следующие обозначения:

S – исследуемая сложная система;

X – имеющиеся статистические данные об S , причем $X \in \mathcal{X}$, где \mathcal{X} – множество всех возможных данных об S ;

Y – результаты применения методов САД к X ;

$d(S)$ – выводы или решения относительно S .

Тогда САД можно представить в виде схемы на рисунке 1:

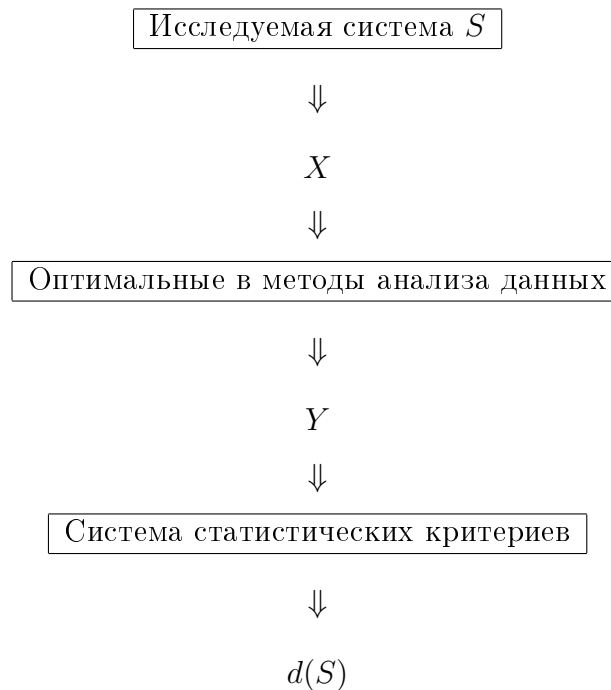


Рисунок 1. Схема статистического анализа данных.

В рамках статистического анализа данных предполагается, что данные о системе S имеют стохастическую природу, а для их описания и анализа используются вероятностно-статистические модели и методы.

В основе САД лежат следующие принципы:

- переход от анализа всех мыслимых значений случайной величины ξ , описывающей функционирование и свойства S (введенное выше множество \mathcal{X}), к анализу случайной выборки X значений из распределения вероятностей ξ ;
- замена случайной выборки X некоторой совокупностью Y вычисленных на ее основе выборочных характеристик (функций от выборки), называемых *статистиками*;
- принятие решений (или статистических выводов) на основе анализа статистик с помощью статистических критериев, справедливых с некоторой заданной вероятностью ошибки.

Статистический анализ данных можно условно разделить на следующие этапы:

Этап 1. Исследование сложной системы и сбор априорной информации.

На этом этапе выявляются основные факторы, определяющие функционирование системы, а также изучаются ее свойства.

Этап 2. Формирование предположений относительно системы и построение ее математической модели.

С учетом результатов, полученных на первом этапе, строится математическая модель системы S и проверяется адекватность этой модели. При этом применяются методы теории вероятностей и математической статистики.

Этап 3. Планирование и проведение экспериментов с системой и сбор статистических данных.

Разрабатывается план экспериментов, в соответствии с которым проводятся эксперименты над математической моделью, полученной на втором этапе. Целью этого этапа является получение статистических данных.

Этап 4. Первичный (разведочный) анализ и сжатие данных.

Проводится предварительный анализ и, если это необходимо и возможно, сжатие данных.

Этап 5. Подбор методов и непосредственное проведение статистического анализа.

Этап 6. Анализ полученных результатов и проверка их адекватности.

Этап 7. Принятие решения об удовлетворительности полученных результатов или необходимости возврата на предыдущие этапы.

Глубина и достоверность статистических выводов и решений зависит от математической модели системы S и методов, применяемых при исследовании этой системы, а также от того, насколько корректно эти методы и модель применяются.

В различных приложениях статистического анализа данных существуют три основных типа данных:

- *пространственные данные* – это совокупность значений $x_1, x_2, \dots, x_n \in \mathbb{R}^N$ анализируемых переменных, полученных для некоторого множества из n ($n > 1$) объектов исследования в фиксированный момент (период) времени.

Пространственные данные характеризуют зависимости между объектами исследования, описываемые выбранными переменными.

- *временной ряд* – это последовательность значений $x_1, x_2, \dots, x_T \in \mathbb{R}^N$ анализируемых переменных, которые соответствуют T ($T > 1$) последовательным моментам (периодам) времени.

Временные ряды характеризуют динамику изменения анализируемых переменных во времени.

- *панельные данные* – это совокупность значений $x_{1t}, x_{2t}, \dots, x_{nt} \in \mathbb{R}^N$ анализируемой переменной, полученных для некоторого множества из n ($n > 1$) объектов исследования в последовательные моменты (периоды) времени $t = 1, 2, \dots, T$.

Панельные данные описывают динамику изменения показателей, характеризующих состояние некоторой группы исследуемых объектов с учетом взаимосвязей между ними.

По характеру компонент (признаков), из которых образованы векторы-наблюдения $x_i = (x_{i1}, \dots, x_{iN})^T$, $i = 1, \dots, n$, выделяют следующие модели данных[11]:

- *непрерывные* – это данные, когда все N признаков имеют совместное абсолютно-непрерывное распределение вероятностей с некоторой плотностью $p(z)$, $z \in R^N$, которая однозначно определяет их совместную функцию распределения $F(z)$, $z \in R^N$;

- *дискретные* – данные, когда наблюдения являются дискретными случайными векторами;

- *комбинированные* – данные, когда часть признаков являются непрерывными случайными величинами, а остальные – дискретными.

Статистический анализ данных обычно проводится с использованием компьютеров и различных пакетов прикладных программ [?], в основу которых положены теоретически обоснованные статистические методы.

1 Предварительный (первичный) статистический анализ данных

Предварительный анализ данных проводится с целью выявления особенностей и свойств исследуемой системы, определения структуры данных и формулировки предположений относительно их модели.

1.1 Математические модели пространственных данных

Случайная выборка из распределения вероятностей

Основные методы САД ориентированы на пространственные данные, адекватной математической моделью которых является случайная выборка.

Пусть состояние исследуемой системы характеризуется N -мерным вектором показателей $\xi \in \mathbb{R}^N$, $N \geq 1$. Предположим, что x_1, \dots, x_n — совокупность (выборка) наблюдаемых значений $\xi \in \mathbb{R}^N$, полученных в экспериментах по наблюдению за n объектами.

Математической моделью i -го наблюдения является случайная величина $x_i = (x_{i1}, \dots, x_{iN})^T \in \mathbb{R}^N$ (скалярная при $N = 1$ или векторная при $N > 1$) с некоторой функцией распределения

$$F_i(z_i) = F_i(z_{i1}, \dots, z_{iN}) = \mathbf{P}\{x_{i1} < z_{i1}, \dots, x_{iN} < z_{iN}\} \equiv \mathbf{P}\{x_i^T < z_i^T\}, \quad (1.1)$$

где $z_i = (z_{i1}, \dots, z_{iN})^T \in \mathbb{R}^N$, $i = 1, \dots, n$, $N \geq 1$. В (1.1) и далее, если не оговорено особо, при $N > 1$ используется формальное обозначение:

$$\mathbf{P}\{x_i^T < z_i^T\} ::= \mathbf{P}\{x_{i1} < z_{i1}, \dots, x_{iN} < z_{iN}\}, \quad i = 1, \dots, n. \quad (1.2)$$

Математической моделью данных, представляющих собой выборку x_1, \dots, x_n , является случайный вектор $X \in \mathbb{R}^{Nn}$ с Nn -мерной функцией распределения

$$\begin{aligned} F_X(z_1, \dots, z_n) &::= F_X(z_{11}, \dots, z_{1N}; z_{21}, \dots, z_{2N}; \dots; z_{n1}, \dots, z_{nN}) = \\ &= \mathbf{P}\{x_1^T < z_1^T, \dots, x_n^T < z_n^T\}. \end{aligned} \quad (1.3)$$

Определение 1.1. Данные $X = (x_1^T, \dots, x_n^T)^T \in \mathbb{R}^{Nn}$ (эквивалентная форма записи $X = \{x_1, \dots, x_n\}$) принято называть случайной выборкой объема n из N -мерного распределения вероятностей $\mathcal{L}\{\xi\}$ с функцией распределения $F(z)$ (или просто — случайной выборкой из распределения $\mathcal{L}\{\xi\}$), если x_1, \dots, x_n являются независимыми одинаково распределенными случайными величинами (н. о. р. с. в.), то есть если для них выполняются предположения:

P_1 . Случайные величины x_1, \dots, x_n одинаково распределены:

$$F_{x_1}(z_1) = \dots = F_{x_n}(z_n) \equiv F(z), \quad z \in \mathbb{R}^N;$$

P_2 . Случайные величины x_1, \dots, x_n независимы в совокупности, при этом:

$$F_X(z_1, \dots, z_n) = \prod_{i=1}^n F_{x_i}(z_i) = (F(z))^n, \quad z \in \mathbb{R}^N.$$

Математические модели данных, альтернативные случайной выборке

Во многих практических задачах предположения P_1 и P_2 о данных, сформулированные в определении 1.1, могут нарушаться. Это приводит к моделям, альтернативным случайным выборкам. В связи с тем, что основная часть методов САД ориентирована на модели данных типа случайная выборка, на этапе предварительного анализа данных необходимо выявить факторы, приводящие к нарушениям предположений P_1 и P_2 . В тех случаях, когда эти факторы могут быть выявлены и устранены, такие модели могут быть представлены в виде случайной выборки из некоторого распределения вероятностей. Для анализа полученных данных могут применяться классические методы САД.

При анализе пространственных данных встречаются следующие альтернативные случайной выборке модели данных:

- преднамеренная (неслучайная) выборка;
- неоднородная выборка из смеси распределений;
- случайная выборка с засорениями.

Преднамеренная выборка – выборка, полученная в результате пристрастного отбора данных.

Такие выборки, очевидно, нельзя считать случайными, поскольку при их формировании нарушаются принципы случайности и независимости экспериментов, необходимых для получения статистических данных. Кроме того, преднамеренные выборки могут быть следствием и непреднамеренных действий. Например, при недостаточной точности измерительного инструмента.

Неоднородная выборка – случайная выборка из смеси нескольких распределений вероятностей.

Неоднородная выборка может быть следствием включения в одну группу объектов, которые существенно отличаются своими функциональными характеристиками. Например: выборка значений показателей единиц продукции, произведенной в различных условиях; выборка значений показателей физического состояния пациентов до и после лечения и т. д.

Такая выборка может быть описана следующей моделью. Пусть в пространстве \mathbb{R}^N , $N \geq 1$, с априорными вероятностями π_0, \dots, π_{L-1} ,

$$\pi_0 + \dots + \pi_{L-1} = 1,$$

наблюдаются объекты из L классов, $L \geq 2$. Случайный вектор наблюдения $x \in \mathbb{R}^N$ для объектов из класса l имеет условную плотность распределения $p_l(z)$, $l = 0, \dots, L-1$. Тогда выборочные наблюдения x_1, \dots, x_n можно рассматривать как случайную выборку из некоторого распределения, которое является смесью распределений и имеет плотность распределения вероятностей вида:

$$p(z_i) = \sum_{l=0}^{L-1} \pi_l p_l(z_i), \quad i = 1, \dots, n. \quad (1.4)$$

Случайная выборка с засорениями – частный случай неоднородной выборки.

Наиболее известная модель данного класса – модель засорений Тьюки. Она может быть описана формулой (1.4) при $L = 2$, $\pi_0 = 1 - \varepsilon$, $\pi_1 = \varepsilon$, где ε – близкий к нулю уровень засорения (доля аномальных наблюдений) $0 < \varepsilon < 0,5$. Выборку с моделью засорений Тьюки можно рассматривать как случайную выборку из распределения вероятностей с плотностью

$$p(z_i) = (1 - \varepsilon)p_0(z_i) + \varepsilon p_1(z_i),$$

где $p_0(z_i)$ – плотность распределения гипотетического распределения основной массы наблюдений, а $p_1(z_i)$ – плотность «засоряющего» распределения, $i = 1, \dots, n$.

1.2 Математические модели статистических наблюдений в динамике

Во многих задачах САД существенным является порядок поступления наблюдений за исследуемой системой. Математическими моделями таких данных являются случайные функции (случайные процессы, временные ряды).

Определение 1.2. *Случайной функцией (случайным процессом) называется параметрическое семейство случайных векторов:*

$$\xi = \xi(t) = \xi(\omega, t) \in \mathbb{R}^N,$$

определенных на одном и том же вероятностном пространстве $(\Omega, \mathcal{F}, \mathbf{P})$, $\omega \in \Omega$, где t – параметр, изменяющийся на некотором множестве $\mathcal{T}: t \in \mathcal{T} \subseteq \mathbb{R}^m$.

В зависимости от значений N и m размерностей пространства наблюдений \mathbb{R}^N и пространства значений параметра \mathbb{R}^m принята классификация случайных функций, приведенная в таблице 1 [12].

Таблица 1 – Классификация случайных функций

Размерность пространства наблюдений, N	Размерность пространства значений параметра, m	
	$m = 1$	$m > 1$
$N = 1$	<i>Случайный процесс</i> (Цена акций [3]; показатель солнечной активности [2])	<i>Случайное поле</i> (Экран монитора в черно-белом режиме, $m = 2$)
$N > 1$	<i>Векторный случайный процесс</i> (Несколько отведений электрокардиограммы, обычно $N = 8$)	<i>Векторное случайное поле</i> (Радиационная карта местности с несколькими радиоизотопами, $m = 2$)

Если множество \mathcal{T} значений параметра t дискретно (конечно или счетно), то такой случайный процесс обычно называют *временным рядом* или *случайной последовательностью*, а одномерный параметр понимают как *время*. Векторный случайный процесс при этом называют *векторным временным рядом*.

2 Функциональные и числовые характеристики вероятностных моделей данных

2.1 Функциональные характеристики вероятностных моделей данных

На этапе разведочного анализа, как правило, исследуются одномерные (маргинальные) распределения вероятностей каждого признака по отдельности. Поэтому, вначале предположим, что наблюдается некоторый одномерный признак (случайная величина) ξ .

Вероятностный закон распределения $\mathcal{L}\{\xi\}$ одномерной случайной величины $\xi \in \mathbb{R}$ описывается следующими функциональными характеристиками [12].

Функция распределения вероятностей:

$$F_{\xi}(z) = \mathbf{P} \{ \xi < z \}, \quad z \in \mathbb{R}. \quad (2.1)$$

Если ξ – дискретная случайная величина, принимающая значения из множества $C = \{c_1, \dots, c_K\}$ с вероятностями

$$p_k = \mathbf{P} \{ x = c_k \}, \quad k = 1, \dots, K, \quad K \leq +\infty,$$

то

$$F_{\xi}(z) = \sum_{k: c_k < z} \mathbf{P} \{ \xi = c_k \} = \sum_{k: c_k < z} p_k, \quad z \in \mathbb{R}. \quad (2.2)$$

Если ξ – непрерывная случайная величина, то

$$F_{\xi}(z) = \int_{-\infty}^z p_{\xi}(y) dy, \quad z \in \mathbb{R}, \quad (2.3)$$

где $p_{\xi}(y)$ – плотность распределения вероятности ξ , $y \in \mathbb{R}$.

Распределение вероятностей дискретной случайной величины:

$$p_k = \mathbf{P} \{ \xi = c_k \}, \quad k = 1, \dots, K, \quad K \leq +\infty, \quad (2.4)$$

где $C = \{c_1, \dots, c_K\}$ – конечное ($K < +\infty$) или счетное ($K = +\infty$) множество возможных значений ξ и

$$\sum_{k=1}^K p_k = 1.$$

Плотность распределения вероятностей непрерывной случайной величины:

$$p_{\xi}(z) = F'_{\xi}(z), \quad z \in \mathbb{R}, \quad (2.5)$$

где $F_{\xi}(z)$ – функция распределения ξ , $z \in \mathbb{R}$.

Характеристическая функция:

$$f_{\xi}(t) = \mathbf{E} \{ e^{it\xi} \} = \int_{-\infty}^{+\infty} e^{itz} dF_{\xi}(z), \quad t \in \mathbb{R}, \quad f_{\xi}(t) \in \mathbb{C}. \quad (2.6)$$

Если ξ – дискретная случайная величина, принимающая значения из множества $C = \{c_1, \dots, c_K\}$ с вероятностями

$$p_k = \mathbf{P} \{ \xi = c_k \}, \quad k = 1, \dots, K, \quad K \leq +\infty,$$

то

$$f_{\xi}(t) = \sum_{k=1}^K e^{itc_k} \cdot p_k, \quad t \in \mathbb{R}. \quad (2.7)$$

Если ξ – непрерывная случайная величина, то

$$f_{\xi}(t) = \int_{-\infty}^{+\infty} e^{ity} p_x(y) dy, \quad t \in \mathbb{R}, \quad (2.8)$$

где $p_{\xi}(y)$ – плотность распределения вероятностей ξ , $y \in \mathbb{R}$.

Если в эксперименте наблюдается N -мерный случайный вектор

$$\xi = (\xi_1, \xi_2, \dots, \xi_N)^{\mathbf{T}} \in \mathbb{R}^N, \quad N \geq 2,$$

то закон распределения вероятностей $\mathcal{L}(\xi)$ описывается следующими функциональными характеристиками.

N -мерная функция распределения:

$$F_{\xi}(z) = F_{\xi_1, \dots, \xi_N}(z_1, \dots, z_N) = \mathbf{P} \{ \xi_1 < z_1, \dots, \xi_N < z_N \}, \quad (2.9)$$

где $\xi, z \in \mathbb{R}^N$.

Распределение вероятностей дискретного случайного N -вектора:

$$p_{k_1, \dots, k_N} = \mathbf{P} \{ \xi = c_{k_1, \dots, k_N} \} = \mathbf{P} \left\{ \xi_1 = c_1^{(k_1)}, \dots, \xi_N = c_N^{(k_N)} \right\}, \quad (2.10)$$

$$k_j = 1, 2, \dots, K_j, \quad j = 1, 2, \dots, N,$$

где конечное или счетное множество возможных значений случайного N -вектора ξ определено соотношениями:

$$C = \{ c_{k_1, \dots, k_N} : k_j = 1, 2, \dots, K_j, j = 1, 2, \dots, N \},$$

$$c_{k_1, \dots, k_N} = \left(c_1^{(k_1)}, \dots, c_N^{(k_N)} \right)^{\mathbf{T}} \in \mathbb{R}^N, \quad k_j = 1, 2, \dots, K_j, j = 1, 2, \dots, N,$$

$K_1 \leq \infty, \dots, K_N \leq \infty$, а

$$\sum_{k_1=1}^{K_1} \dots \sum_{k_N=1}^{K_N} p_{k_1, \dots, k_N} = 1.$$

N -мерная плотность распределения вероятностей:

$$p_{\xi}(z) = p_{\xi}(z_1, \dots, z_N) = \frac{\partial^N F_{\xi}(z_1, \dots, z_N)}{\partial z_1 \dots \partial z_N},$$

где $F_{\xi}(z)$ – N -мерная функция распределения случайного N -вектора ξ , $z \in \mathbb{R}^N$.

Характеристическая функция N -мерного распределения вероятностей:

$$f_{\xi}(t) = \mathbf{E} \left\{ e^{it^{\mathbf{T}} \xi} \right\} = \mathbf{E} \left\{ e^{i \langle t, \xi \rangle} \right\} = \mathbf{E} \exp \left\{ i \sum_{j=1}^N t_j \xi_j \right\} =$$

$$= \int_{\mathbb{R}^N} \exp \left\{ i \sum_{j=1}^N t_j z_j \right\} dF_{\xi}(z), \quad t, z \in \mathbb{R}^N, \quad f_{\xi}(t) \in \mathbb{C}, \quad (2.11)$$

где $\langle t, \xi \rangle$ – скалярное произведение векторов t и ξ .

2.2 Числовые характеристики вероятностных моделей данных

Числовые характеристики распределения вероятностей $\mathcal{L}\{\xi\}$ одномерной случайной величины $\xi \in \mathbb{R}$ позволяют составить наглядное представление об этом распределении. При разведочном анализе данных обычно вычисляют следующие числовые характеристики [12].

Характеристики положения

Математическое ожидание:

$$\mu = \mathbf{E}\{\xi\}.$$

Если ξ – дискретная случайная величина, принимающая значения из множества $C = \{c_1, \dots, c_K\}$ с вероятностями

$$p_k = \mathbf{P}\{\xi = c_k\}, \quad k = 1, \dots, K, \quad K \leq +\infty,$$

то

$$\mu = \sum_{k=1}^K c_k \cdot p_k. \quad (2.12)$$

Если ξ – непрерывная случайная величина, то

$$\mu = \int_{-\infty}^{+\infty} y p_\xi(y) dy, \quad (2.13)$$

где $p_\xi(y)$ – плотность распределения вероятности ξ ($y \in \mathbb{R}$).

Наибольшее и наименьшее значения случайной величины ξ :

$$c_+ = \max\{\xi\}, \quad c_- = \min\{\xi\}. \quad (2.14)$$

Если ξ – дискретная случайная величина, принимающая значения из множества $C = \{c_1, \dots, c_K\}$

$$c_+ = \max\{c_1, \dots, c_K\}, \quad c_- = \min\{c_1, \dots, c_K\}.$$

Если ξ – непрерывная случайная величина, то

$$c_+ = \max\{z : p_\xi(z) \neq 0\}, \quad c_- = \min\{z : p_\xi(z) \neq 0\},$$

где $p_\xi(z)$ – плотность распределения вероятности ξ , $y \in \mathbb{R}$.

Мода распределения вероятностей $\mathcal{L}\{\xi\}$:

Если ξ – дискретная случайная величина, принимающая значения из множества $C = \{c_1, \dots, c_K\}$ с вероятностями

$$p_k = \mathbf{P}\{x = c_k\}, \quad k = 1, \dots, K, \quad K \leq +\infty,$$

то

$$M = \text{mod } \xi = c_{k^*}, \quad \text{где } k^* = \arg \max\{p_1, \dots, p_K\}. \quad (2.15)$$

Если ξ – непрерывная случайная величина с плотностью распределения вероятности $p_\xi(y)$, $y \in \mathbb{R}$, то

$$M = \text{mod } \xi = \arg \max_{y \in \mathbb{R}} p_\xi(y). \quad (2.16)$$

Медиана распределения вероятностей $\mathcal{L}\{\xi\}$:

$$m = \text{med } \xi = z_{0,5} = F_\xi^{-1}(0,5), \quad (2.17)$$

где $F_\xi(z)$ – функция распределения случайной величины ξ .

Характеристики рассеяния

Дисперсия:

$$\sigma^2 = \mathbf{D}\{\xi\} = \mathbf{E}\{(\xi - \mu)^2\}. \quad (2.18)$$

Среднее квадратическое отклонение:

$$\sigma = \sqrt{\mathbf{D}\{\xi\}}. \quad (2.19)$$

Размах значений случайной величины ξ :

$$R = c_+ - c_-, \quad (2.20)$$

где c_+ и c_- определены в (2.14).

Интерквартильный размах:

$$IR = z_{\frac{3}{4}} - z_{\frac{1}{4}}, \quad (2.21)$$

где $z_{\frac{k}{4}}$ — k -я квартиль распределения вероятностей $\mathcal{L}\{\xi\}$, то есть

$$z_{\frac{k}{4}} = F_{\xi}^{-1}\left(\frac{k}{4}\right), \quad k \in \{1, 3\}. \quad (2.22)$$

Характеристики формы

Асимметрия распределения вероятностей $\mathcal{L}\{\xi\}$:

$$\beta_1 = \frac{\alpha_3}{\sigma^3}, \quad (2.23)$$

где α_3 — начальный момент 3-го порядка случайной величины ξ

$$\alpha_3 = \mathbf{E}\{\xi^3\},$$

а $\sigma = \sqrt{\mathbf{D}\{\xi\}}$.

Экссесс распределения вероятностей $\mathcal{L}\{\xi\}$:

$$\beta_2 = \frac{\mu_4}{\sigma^4} - 3, \quad (2.24)$$

где μ_4 — центральный момент 4-го порядка случайной величины ξ

$$\mu_4 = \mathbf{E}\{(\xi - \mu)^4\},$$

а $\sigma = \sqrt{\mathbf{D}\{\xi\}}$.

Если наблюдаемая в эксперименте случайная величина ξ представляет собой случайный вектор размерности N , то есть

$$\xi = (\xi_1, \dots, \xi_N)^{\mathbf{T}} \in \mathbb{R}^N,$$

то числовыми характеристиками этого случайного вектора, которые описывают N -мерный закон распределения $\mathcal{L}\{\xi\}$ являются

Вектор математического ожидания:

$$m = (m_1, \dots, m_N)^{\mathbf{T}} \in \mathbb{R}^N, \quad (2.25)$$

где

$$m_j = \mathbf{E}\{\xi_j\}, \quad j = 1, \dots, N. \quad (2.26)$$

Центральный смешанный момент порядка r :

$$\mu_{r_1, \dots, r_N} = \mathbf{E}\{(\xi_1 - m_1)^{r_1} \cdot \dots \cdot (\xi_N - m_N)^{r_N}\}, \quad (2.27)$$

где r_1, \dots, r_N — целые неотрицательные целые числа, задающие разбиение числа $r \in \mathbb{N}$, то есть

$$r = r_1 + \dots + r_N.$$

Ковариационная матрица:

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1N} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{N1} & \sigma_{N2} & \dots & \sigma_{NN} \end{pmatrix} \quad (2.28)$$

где σ_{jk} — ковариация случайных величин ξ_j и ξ_k (центральный смешанный момент второго порядка):

$$\sigma_{jk} = \mathbf{Cov}\{\xi_j, \xi_k\} = \mathbf{E}\{(\xi_j - m_j) \cdot (\xi_k - m_k)\}, \quad j, k = 1, \dots, N. \quad (2.29)$$

Очевидно, что σ_{jj} — дисперсия случайной величины ξ_j , $j = 1, \dots, N$:

$$\sigma_{jj} = \sigma_j^2.$$

Корреляционная матрица:

$$\mathcal{P} = \begin{pmatrix} \rho_{11} & \rho_{12} & \dots & \rho_{1N} \\ \rho_{21} & \rho_{22} & \dots & \rho_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{N1} & \rho_{N2} & \dots & \rho_{NN} \end{pmatrix} \quad (2.30)$$

где ρ_{jk} — коэффициент корреляции случайных величин ξ_j и ξ_k :

$$\rho_{jk} = \frac{\sigma_{jk}}{\sigma_{jj} \cdot \sigma_{kk}}, \quad j, k = 1, \dots, N. \quad (2.31)$$

Очевидно, что $\rho_{jj} = 1$, $j = 1, \dots, N$.

3 Статистическое оценивание функциональных и числовых характеристик вероятностных моделей данных

3.1 Основные понятия теории статистического оценивания параметров

Пусть в результате проведения экспериментов наблюдается случайная выборка $X = \{x_1, \dots, x_n\}$ объема n из некоторого N -мерного распределения вероятностей $\mathcal{L}\{\xi\}$, $\xi \in \mathbb{R}^N$, информация о котором отсутствует полностью или частично. Основной задачей предварительного анализа является восстановление этого распределения по выборке случайных данных. Известны два основных подхода к решению данной задачи: параметрический и непараметрический.

Параметрический подход применяется тогда, когда распределение вероятностей $\mathcal{L}\{\xi\}$, $\xi \in \mathbb{R}^N$, известно с точностью до параметров, то есть известно с точностью до параметров аналитическое представление функции распределения

$$F_\xi(z) = F_0(z, \theta),$$

или других функциональных характеристик распределения вероятностей $\mathcal{L}\{\xi\}$, где $F_0(\cdot)$ – истинная (модельная) функция распределения, а $\theta = (\theta_1, \dots, \theta_m)^T \in \mathbb{R}^m$ – вектор неизвестных параметров, $m \geq 1$.

В этом случае по выборке $X = \{x_1, \dots, x_n\}$ методами математической статистики (методом максимального правдоподобия, методом моментов, методом наименьших квадратов и т. д. [12]) вычисляются оценки $\hat{\theta}_1, \dots, \hat{\theta}_m$ неизвестных параметров $\theta_1, \dots, \theta_m$. Далее, в соответствии с так называемым «подстановочным принципом», полученные оценки $\hat{\theta}_1, \dots, \hat{\theta}_m$ подставляются в соответствующее аналитическое выражение для функции распределения $F_0(z, \theta)$ (или другой известной функциональной характеристики), что приводит к ее статистической оценке $\hat{F}(z)$:

$$\hat{F}(z) \equiv \hat{F}_\xi(z) = F_0(z, \hat{\theta}) = F_0(z; \hat{\theta}_1, \dots, \hat{\theta}_m) \quad z \in \mathbb{R}^N.$$

Замечание 3.0.1. Теоретической основой параметрического метода, а также и подстановочного принципа, являются статистические свойства оценок и предельные теоремы теории вероятностей [12].

Непараметрический подход применяется тогда, когда распределение вероятностей $\mathcal{L}\{\xi\}$, $\xi \in \mathbb{R}^N$, неизвестно, то есть неизвестны его функциональные характеристики. В этом случае по выборке $X = \{x_1, \dots, x_n\}$ вычисляют соответствующие выборочные функциональные характеристики.

При вычислении статистических оценок неизвестных функциональных и числовых характеристик распределения вероятностей $\mathcal{L}\{\xi\}$ по случайной выборке $X = \{x_1, \dots, x_n\}$ объема n всегда возникает вопрос о требованиях, которые следует предъявить к этим статистикам, чтобы они были в каком-то определенном смысле близкими к истинному значению. Эти требования формулируются обычно с помощью следующих статистических свойств оценок: состоятельности, несмещенности и эффективности [12].

Пусть $\theta_n = \hat{\theta}(X)$ – оценка неизвестного параметра θ , вычисленная по случайной выборке $X = \{x_1, \dots, x_n\}$ объема n из некоторого распределения вероятностей $\mathcal{L}\{\xi\}$. Тогда θ_n называется:

– *состоятельной* оценкой параметра θ , если

$$\theta_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \theta;$$

– строго состоятельной оценкой параметра θ , если

$$\theta_n \xrightarrow[n \rightarrow \infty]{\mathbf{P}=1} \theta;$$

– несмещенной оценкой параметра θ , если для любого $n \geq 1$

$$\mathbf{E} \{\theta_n\} = \theta;$$

– асимптотически несмещенной оценкой параметра θ , если

$$\lim_{n \rightarrow \infty} \mathbf{E} \{\theta_n - \theta\} = 0;$$

– эффективной оценкой параметра θ , если θ_n – несмещенная оценка и для любых $\theta \in \mathbb{R}^m$ и $n \geq 1$ эффективность \mathcal{E}_n равна единице:

$$\mathcal{E}_n = \frac{1}{|V||\mathcal{I}_n|} = 1,$$

где V – матрица вариаций оценки θ_n [12], а \mathcal{I}_n – информационная матрица Фишера [12] для всей выборки X .

3.2 Выборочные функциональные характеристики и их статистические свойства

Пусть $X = \{x_1, \dots, x_n\}$ – случайная выборка объема n из некоторого одномерного распределения вероятностей $\mathcal{L}\{\xi\}$. Рассмотрим задачу построения по выборке X статистических оценок функциональных характеристик этого распределения: функции распределения $F(\cdot)$, распределения вероятностей $\{p_k, k = 1, \dots, K\}$, плотности распределения вероятностей $p(\cdot)$, характеристической функции $f(\cdot)$.

Вариационный ряд выборки и порядковые статистики

Определение 3.1. Вариационным рядом выборки $X = \{x_1, \dots, x_n\}$ объема n , называется последовательность $x_{(1)}, x_{(2)}, \dots, x_{(n)}$, упорядоченная по неубыванию, то есть

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}, \quad (3.1)$$

где $x_{(k)}$ – k -я порядковая статистика [12], $k \in \{1, \dots, n\}$.

Вариационный ряд является одним из стандартных способов представления случайной выборки.

Выборочная функция распределения

Определение 3.2. Выборочной функцией распределения или эмпирической функцией распределения, построенной по выборке X объема n , называется статистика

$$\hat{F}(z) ::= F_n(z) = \frac{1}{n} \sum_{i=1}^n \mathbf{I}(z - x_i) = \frac{L_n(z)}{n}, \quad z \in \mathbb{R}, \quad (3.2)$$

где $L_n(z)$ – число выборочных значений из множества $\{x_1, \dots, x_n\}$, для которых $x_i < z$, $i = 1, \dots, n$:

$$L_n(z) = \sum_{i=1}^n \mathbf{I}(z - x_i),$$

а $\mathbf{I}(z)$ – единичная функция Хевисайда:

$$\mathbf{I}(z) = \begin{cases} 1 & \text{при } z > 0, \\ 0 & \text{при } z \leq 0. \end{cases} \quad (3.3)$$

Для построения графика выборочной функции распределения используют вариационный ряд (3.1) выборки X . Схематический график выборочной функции распределения представлен на рисунке 2 [12].

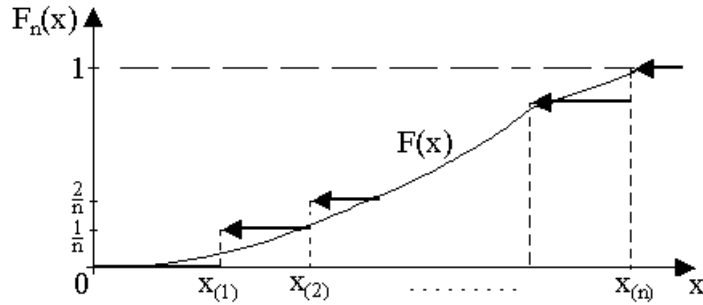


Рисунок 2. Выборочная функция распределения

Относительные частоты

Пусть $X = \{x_1, \dots, x_n\}$ – случайная выборка объема n из некоторого одномерного дискретного распределения вероятностей $\mathcal{L}\{\xi\}$, а наблюдаемая в эксперименте случайная величина ξ принимает значения из множества $C = \{c_1, \dots, c_K\}$, $K \leq \infty$. Тогда, согласно закону больших чисел [12], в качестве оценок $\{\hat{p}_k, k = 1, \dots, K\}$ для неизвестных вероятностей

$$p_k = \mathbf{P}\{x = c_k\}, \quad k = 1, \dots, K,$$

рассматривают относительные частоты

$$\hat{p}_k = \nu_k = \frac{1}{n} \sum_{i=1}^n \delta_{x_i, c_k}, \quad k = 1, \dots, K, \quad (3.4)$$

где

$$\delta_{a,b} = \begin{cases} 1, & \text{если } a = b, \\ 0, & \text{если } a \neq b. \end{cases} \quad (3.5)$$

Гистограмма

Предположим, что плотность распределения вероятностей $p(z)$, $z \in \mathbb{R}$, наблюдаемой непрерывной случайной величины ξ сосредоточена на отрезке $\Gamma = [x_-, x_+]$, а вне этого отрезка она равна нулю. Зададим натуральное число M и осуществим разбиение Γ на M частей точками деления:

$$x_- = b_0 < b_1 < \dots < b_M = x_+.$$

Обозначим m -ю ячейку этого разбиения через Γ_m :

$$\Gamma_m ::= [b_{m-1}, b_m), \quad m = 1, \dots, M, \quad M \geq 2.$$

В дальнейшем Γ_m будем называть m -й ячейкой гистограммы, $m = 1, \dots, M$.

Замечание 3.2.1. При оценивании плотности распределения $p(z)$ по выборке $X = \{x_1, \dots, x_n\}$ объема n как правило полагают:

$$x_- = \min \{x_1, \dots, x_n\}, \quad x_+ = \max \{x_1, \dots, x_n\},$$

а

$$M ::= M_n = [\log_2 n] + 1.$$

Введем следующие обозначения:

- Δ_m – «размер» (мера Лебега) m -й ячейки гистограммы:

$$\Delta_m ::= \mu(\Gamma_m) = b_m - b_{m-1}, \quad m = 1, \dots, M;$$

- ν_m – число выборочных значений, попавших в m -ю ячейку:

$$\nu_m = \sum_{i=1}^n \mathbf{I}_{\Gamma_m}(x_i), \quad m = 1, \dots, M,$$

где $\mathbf{I}_A(z)$ – индикатор множества A :

$$\mathbf{I}_A(z) = \begin{cases} 1, & \text{если } z \in A, \\ 0, & \text{если } z \notin A. \end{cases} \quad (3.6)$$

Определение 3.3. Гистограммой или гистограммной оценкой плотности распределения $p(z)$ называется статистика

$$\hat{p}(z) = \sum_{m=1}^M \frac{\nu_m}{n\Delta_m} \mathbf{I}_{\Gamma_m}(z), \quad z \in \mathbb{R}.$$

График гистограммы с наложенным графиком истинной плотности распределения приведен на рисунке 3.

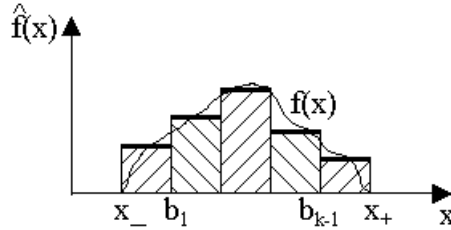


Рисунок 3. Гистограмма

Заметим, что гистограмма, как и плотность распределения, удовлетворяет условию нормировки:

$$\int_{-\infty}^{\infty} \hat{p}(z) dz = \int_{-\infty}^{\infty} p(z) dz = 1.$$

Отметим также, что гистограмма – смещенная и несостоятельная оценка плотности распределения [12]. Для ее состоятельности необходимо, чтобы при $n \rightarrow +\infty$ ячейки «измельчались» специальным образом [12]:

$$M = M_n \xrightarrow{n \rightarrow \infty} +\infty, \quad \max_{1 \leq m \leq M_n} \Delta_m \xrightarrow{n \rightarrow \infty} 0.$$

Выборочная характеристическая функция

Определение 3.4. Выборочной характеристической функцией, построенной по выборке X объема n , называется комплекснозначная функция действительной переменной t :

$$\hat{f}(t) ::= f_n(t) = \frac{1}{n} \sum_{k=1}^n e^{itx_k}, \quad t \in \mathbb{R}. \quad (3.7)$$

3.3 Выборочные числовые характеристики и их статистические свойства

Пусть $X = \{x_1, \dots, x_n\}$ – выборка объема n из некоторого одномерного распределения вероятностей $\mathcal{L}\{\xi\}$, $F_\xi(z)$ и $F_n(z)$ – соответственно теоретическая и эмпирическая функции распределения. $F_n(z)$ можно рассматривать как функцию распределения некоторой дискретной случайной величины x , принимающей n значений $\{x_1, \dots, x_n\}$ с вероятностями, равными $\frac{1}{n}$, причем, если какое-либо из этих значений встретится k раз в выборке X , то этому значению соответствует вероятность $\frac{k}{n}$ [7]. Так же как для теоретического распределения $\mathcal{L}\{\xi\}$, так и для эмпирического распределения, обозначим его $\mathcal{L}\{x\}$, связанного с выборкой X , вводятся аналогичные числовые характеристики, которые принято называть выборочными характеристиками.

Эмпирическая функция распределения $F_n(z)$ является сильно состоятельной несмещенной оценкой для теоретической функции распределения $F_\xi(z)$ [12]. Следовательно, $F_n(z)$ – статистический аналог $F_\xi(z)$ и выборочная характеристика – статистический аналог соответствующей теоретической характеристики [7]. В общем случае, если $g = \mathbf{E}\{g(\xi)\}$ – некоторая теоретическая характеристика наблюдаемой случайной величины ξ , то ее статистический аналог, то есть соответствующая выборочная характеристика, вычисляется по формуле [7]:

$$\hat{g} ::= g_n = \frac{1}{n} \sum_{i=1}^n g(x_i). \quad (3.8)$$

Таким образом, *выборочные начальные и выборочные центральные моменты* вычисляются, соответственно, по формулам:

$$\hat{\alpha}_r = \frac{1}{n} \sum_{i=1}^n x_i^r, \quad (3.9)$$

$$\hat{\mu}_r = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\alpha}_1)^r. \quad (3.10)$$

Причем выборочный начальный момент r -го порядка $\hat{\alpha}_r$ является сильно состоятельной оценкой для начального момента r -го порядка α_r [12].

На этапе предварительного анализа по случайной выборке $X = \{x_1, \dots, x_n\}$ объема n обычно вычисляют следующие выборочные числовые характеристики:

Характеристики положения

Минимальное и максимальное значения:

$$\hat{c}_- = x_{(1)} = \min_{1 \leq i \leq n} x_i, \quad \hat{c}_+ = x_{(n)} = \max_{1 \leq i \leq n} x_i. \quad (3.11)$$

Выборочное среднее:

$$\bar{x} = \hat{\nu}_1 = \frac{1}{n} \sum_{i=1}^n x_i. \quad (3.12)$$

\bar{x} ($\hat{\nu}_1$) – статистическая оценка математического ожидания $\mu = \mathbf{E}\{\xi\}$.

Выборочная мода:

$$\hat{M} = x_{k^*}, \text{ где } k^* = \arg \max \{\hat{p}_1, \dots, \hat{p}_K\}, \quad (3.13)$$

где \hat{p}_k , $k = 1, \dots, K$, определены в (3.4).

Выборочная медиана:

$$\hat{m} = \begin{cases} x_{(k+1)}, & n = 2k + 1; \\ \frac{x_{(k)} + x_{(k+1)}}{2}, & n = 2k, \end{cases} \quad (3.14)$$

где $x_{(k)}$ – k -я порядковая статистика.

\hat{m} – статистическая оценка медианы m .

Характеристики рассеяния

Выборочная дисперсия:

$$\hat{\sigma}^2 ::= \hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (3.15)$$

Выборочная дисперсия $\hat{\sigma}^2$ сходится по вероятности к теоретической дисперсии σ^2 [7]. Следовательно, $\hat{\sigma}^2$ – состоятельная оценка для σ^2 . Но так как $\mathbf{E} \{\hat{\sigma}^2\} \neq \sigma^2$, то $\hat{\sigma}^2$ – смещенная оценка для σ^2 .

Несмещенной состоятельной оценкой для σ^2 является статистика [7]:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (3.16)$$

При этом s – оценка среднеквадратического отклонения σ .

Размах выборки:

$$\hat{R} = \hat{c}_+ - \hat{c}_-, \quad (3.17)$$

где \hat{c}_+ и \hat{c}_- определены в (3.11).

Интерквартильный размах выборки:

$$I\hat{R} = \hat{z}_{\frac{3}{4}} - \hat{z}_{\frac{1}{4}}, \quad (3.18)$$

где $\hat{z}_{\frac{k}{4}}$ – выборочная k -я квартиль:

$$\hat{z}_{\frac{k}{4}} = F_n^{-1} \left(\frac{k}{4} \right), \quad k \in \{1, 3\}. \quad (3.19)$$

Интервал концентрации:

$$[\bar{x} - 3s, \bar{x} + 3s] \quad \left(\mathbf{P} \{x_i \in [\mu - 3\sigma, \mu + 3\sigma]\} \geq \frac{8}{9} \right).$$

Характеристики формы

Выборочная асимметрия:

$$\hat{\beta}_1 = \frac{\hat{\alpha}_3}{s^3}, \quad (3.20)$$

где $\hat{\alpha}_3$ – выборочный начальный момент 3-го порядка случайной величины, а s – оценка среднеквадратического отклонения.

Выборочный эксцесс:

$$\hat{\beta}_2 = \frac{\hat{\mu}_4}{s^4} - 3, \quad (3.21)$$

где $\hat{\mu}_4$ – выборочный центральный момент 4-го порядка.

3.4 Проблема сжатия данных. Метод главных компонент

В многомерном статистическом анализе исследуемая система описывается вектором, размерность N которого в прикладных задачах может быть достаточно большой. Анализировать такие данные достаточно сложно. Кроме того, те факторы, от которых интересующий исследователя признак не зависит (или зависит в незначительной степени), ухудшают свойства статистических процедур. В частности, включение этих факторов в анализ увеличивает дисперсию оценок параметров и характеристик распределений. Поэтому вполне естественным является желание перейти от многомерной выборки к данным небольшой размерности. Возникает проблема *сжатия данных*.

Одним из наиболее часто используемых методов снижения размерности является метод главных компонент. Основная идея этого метода состоит в последовательном выявлении факторов, в которых данные имеют наибольший разброс.

Пусть в пространстве \mathbb{R}^N зарегистрирована случайная выборка $X = \{x_1, \dots, x_n\}$ объема n : $x_i \in \mathbb{R}^N$, $i = 1, \dots, n$.

Необходимо преобразовать исходную выборку $X = \{x_1, \dots, x_n\}$ в «сжатую» выборку $Y = \{y_1, \dots, y_n\}$: $y_i \in \mathbb{R}^m$, $i = 1, \dots, n$, с меньшим числом признаков m ($m < N$), которые несут информацию об исходных наблюдениях с минимальными (наперед заданными) потерями.

Формально решение задачи сжатия данных сводится к поиску борелевского преобразования:

$$y = B(x) : \mathbb{R}^N \rightarrow \mathbb{R}^m, \quad (3.22)$$

переводящего исходное наблюдение $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$ в «сжатое» наблюдение $y = (\tilde{y}_1, \dots, \tilde{y}_m)^T \in \mathbb{R}^m$. Выборка Y получается из выборки X преобразованием каждого из n наблюдений: $y_i = B(x_i)$, $i = 1, \dots, n$.

Преобразование в (3.22), вообще говоря, нелинейное. Однако на практике обычно используют линейные преобразования:

$$y = Bx : \mathbb{R}^N \rightarrow \mathbb{R}^m,$$

где B – $(m \times N)$ -матрица. Проблема заключается в выборе матрицы B .

Согласно методу главных компонент, случайный N -вектор-наблюдение $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$ с ковариационной $(N \times N)$ -матрицей

$$\Sigma = \mathbf{Cov}\{x, x\} = \mathbf{E}\{(x - \mu)(x - \mu)^T\},$$

$\mu = \mathbf{E}\{x\} \in \mathbb{R}^N$ – вектор математического ожидания, подвергается следующему линейному преобразованию:

$$y = (\tilde{y}_1(x), \dots, \tilde{y}_N(x))^T,$$

которое покомпонентно записывается в виде:

$$\tilde{y}_k = \tilde{y}_k(x) = \Psi_k^T x, \quad k = 1, \dots, N, \quad (3.23)$$

где $\{\Psi_k, k = 1, \dots, N\}$ – ортонормированные собственные векторы ковариационной матрицы Σ , удовлетворяющие соотношениям

$$\Sigma \Psi_k = \lambda_k \Psi_k, \quad k = 1, \dots, N; \quad (3.24)$$

$$\Psi_k^T \Psi_l = \delta_{kl}, \quad k, l = 1, \dots, N, \quad (3.25)$$

и соответствующие упорядоченным по убыванию собственным числам $\{\lambda_k, k = 1, \dots, N\}$ матрицы Σ :

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N.$$

Полученные таким образом случайные величины $\tilde{y}_1, \dots, \tilde{y}_N$ называют *главными компонентами* для исходного наблюдения $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T$. Исследуем вероятностные свойства главных компонент.

Теорема 3.1. *Главные компоненты $\tilde{y}_1, \dots, \tilde{y}_N$, удовлетворяющие (3.23)–(3.25), некоррелированы:*

$$\text{Cov}\{\tilde{y}_k, \tilde{y}_j\} = 0, \quad k \neq j, \quad k, j = 1, \dots, N,$$

а их дисперсии равны соответствующим собственным числам:

$$\mathbf{D}\{\tilde{y}_k\} = \lambda_k, \quad k = 1, \dots, N. \quad (3.26)$$

Доказательство. С учетом (3.23)–(3.25) вычислим ковариацию случайных величин \tilde{y}_k и \tilde{y}_j :

$$\begin{aligned} \text{Cov}\{\tilde{y}_k, \tilde{y}_j\} &= \mathbf{E}\{(\tilde{y}_k - \mathbf{E}\{\tilde{y}_k\})(\tilde{y}_j - \mathbf{E}\{\tilde{y}_j\})\} = \\ &= \Psi_k^T \mathbf{E}\{(x - \mu)(x - \mu)^T\} \Psi_j = \Psi_k^T \Sigma \Psi_j = \lambda_j \Psi_k^T \Psi_j = \delta_{kj} \lambda_j, \quad k, j = 1, \dots, N. \end{aligned}$$

□

Следствие 3.1. *Суммарная дисперсия исходных признаков равна суммарной дисперсии главных компонент:*

$$\sum_{k=1}^N \mathbf{D}\{\tilde{x}_k\} = \sum_{k=1}^N \mathbf{D}\{\tilde{y}_k\} = \sum_{k=1}^N \lambda_k. \quad (3.27)$$

Доказательство. очевидно и следует из (3.26) и известного свойства матрицы [9, 5]:

$$\sum_{k=1}^N \mathbf{D}\{\tilde{x}_k\} = \text{tr}(\Sigma) = \sum_{k=1}^N \lambda_k.$$

□

Соотношения (3.26), (3.27) позволяют предложить критерий выбора *информативных признаков* в пространстве главных компонент, который состоит в том, что признаки, имеющие малые дисперсии, отбрасываются, а рассматриваются лишь m первых главных компонент из $\tilde{y}_1, \dots, \tilde{y}_m, \dots, \tilde{y}_N$ ($\lambda_1 \leq \dots \leq \lambda_m \leq \dots \leq \lambda_N$).

Число m обычно определяется по наперед заданной малой величине $\varepsilon \in [0, 1)$:

$$m = m(\varepsilon) = \min\{k : 1 - \varkappa(k) \leq \varepsilon, \quad k = 1, \dots, N\}; \quad (3.28)$$

$$\varkappa(k) = \frac{\sum_{i=1}^k \lambda_i}{\sum_{j=1}^N \lambda_j},$$

где $0 < \varkappa(k) \leq 1$ в (3.28) – относительная доля суммарной дисперсии первых k главных компонент. Чем ближе $\varkappa(k)$ к единице (а $(1 - \varkappa(k))$ – к нулю), тем меньше потери информации при сжатии данных ($\varkappa(N) = 1$).

Замечание 3.4.1. *На практике обычно истинное значение ковариационной матрицы Σ неизвестно, и вместо нее в методе главных компонент используется оценка по исходной выборке $X = \{x_1, \dots, x_n\}$:*

$$S = \frac{1}{n-1} \sum_{t=1}^n (x_t - \bar{x})(x_t - \bar{x})^T, \quad \bar{x} = \frac{1}{n} \sum_{t=1}^n x_t.$$

Замечание 3.4.2. Метод главных компонент применяется не только для сжатия данных. Он позволяет декоррелировать их – строить выборки из наблюдений с некоррелированными признаками. С его помощью также производят визуализацию многомерных данных: на плоскости ($m = 2$) и в трехмерном пространстве ($m = 3$). При этом учитывается информация, заключенная во всех исходных признаках, чего нельзя сказать о диаграммах рассеяния, на которых исходные признаки отображаются попарно (на плоскости) и тройками (в трехмерном пространстве).

Замечание 3.4.3. Недостаток метода главных компонент состоит в том, что главные компоненты, по сравнению с исходными признаками, не имеют на практике физической интерпретации. С этой точки зрения лучше напрямую выбирать информативные признаки из множества исходных признаков и не производить их функциональных преобразований, то есть осуществлять так называемый «прямой отбор» информативных признаков.

4 Многомерное нормальное (гауссовское) распределение как модель многомерных данных и его свойства

Нормальное распределение играет особую роль в прикладных задачах теории вероятностей и математической статистики в следствие следующих причин:

- разнообразные статистические данные с хорошей степенью точности можно считать выборками из нормального распределения. Примерами могут служить помехи в электроаппаратуре, ошибки измерений, разброс попадания снарядов при стрельбе по заданной цели, рост наудачу взятого человека, скорость реакции на раздражитель и т.д.;
- если на отклонение исследуемой случайной величины ξ от некоторого заданного значения влияет множество различных факторов, причем влияние каждого из них вносит небольшой вклад в это отклонение, а их действия независимы или почти независимы, то можно предполагать, что ξ имеет нормальное распределение;
- в силу центральной предельной теоремы и ее формулировок для частных случаев ([4], [6], [10]) распределение целого ряда широко распространенных в статистике функций от случайных величин (статистик, оценок) хорошо аппроксимируется нормальным распределением.

4.1 Многомерное нормальное распределение и его свойства

Определение 4.1. *Непрерывный случайный N -мерный вектор $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$ с конечными моментами второго порядка $\mathbf{E}\{\tilde{x}_j^2\} < +\infty$, $j = 1, \dots, N$, имеет невырожденное многомерное (N -мерное) нормальное распределение:*

$$\mathcal{L}\{x\} = \mathcal{N}_N(\mu, \Sigma),$$

с N -вектором математического ожидания

$$\mu = \mathbf{E}\{x\} = (\mu_1, \dots, \mu_N)^T \in \mathbb{R}^N$$

и невырожденной $(N \times N)$ -ковариационной матрицей

$$\Sigma = \mathbf{E}\{(x - \mu)(x - \mu)^T\}, \quad |\Sigma| \neq 0,$$

если его плотность распределения вероятностей задается соотношением

$$p(z) = n_N(z|\mu, \Sigma) = (2\pi)^{-\frac{N}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(z - \mu)^T \Sigma^{-1}(z - \mu)\right). \quad (4.1)$$

Замечание 4.1.1. *В случае вырожденной ковариационной матрицы Σ ($|\Sigma| = 0$) иногда также определяют вырожденное многомерное нормальное распределение, которое, однако, не описывается плотностью (6.24). Считается, что случайный N -вектор $x \in \mathbb{R}^N$ имеет вырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| = 0$), если он может быть представлен в виде $x = By + b$, где случайный p -вектор $y \in \mathbb{R}^p$ ($p < N$) имеет невырожденное p -мерное нормальное распределение, $(p \times N)$ -матрица линейного преобразования B – полный ранг: $\text{rank}(B) = p$, а p -вектор сдвига $b = \mu - B\mathbf{E}\{y\}$.*

Многомерное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$ однозначно определяется своими параметрами:

- вектором математического ожидания $\mu = (\mu_1, \dots, \mu_N)^T \in \mathbb{R}^N$, где $\mu_j = \mathbf{E}\{\tilde{x}_j\}$ – математическое ожидание j -й компоненты $j = 1, \dots, N$;

– ковариационной матрицей

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1N} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{N1} & \sigma_{N2} & \dots & \sigma_{NN} \end{pmatrix} \quad (4.2)$$

где σ_{jk} – ковариация случайных величин \tilde{x}_j и \tilde{x}_k :

$$\sigma_{jk} = \mathbf{E}\{(\tilde{x}_j - m_j) \cdot (\tilde{x}_k - m_k)\}, \quad j, k = 1, \dots, N, \quad (4.3)$$

а σ_{jj} – дисперсия j -й компоненты:

$$\sigma_{jj} = \mathbf{D}\{\tilde{x}_j\} = \mathbf{E}\{(\tilde{x}_j - \mu_j)^2\}, \quad j = 1, \dots, N.$$

Все характеристики (функциональные и числовые) многомерного нормального распределения $\mathcal{N}_N(\mu, \Sigma)$ можно определить используя параметры μ и Σ в том числе:

– коэффициент корреляции:

$$\rho_{jk} = \frac{\mathbf{Cov}\{\tilde{x}_j, \tilde{x}_k\}}{\sqrt{\mathbf{D}\{\tilde{x}_j\}\mathbf{D}\{\tilde{x}_k\}}} = \frac{\sigma_{jk}}{\sqrt{\sigma_{jj}\sigma_{kk}}} \in [-1, 1], \quad j, k = 1, \dots, N;$$

– центральные моменты третьего порядка:

$$\mathbf{E}\{(\tilde{x}_j - \mu_j)(\tilde{x}_k - \mu_k)(\tilde{x}_l - \mu_l)\} = 0, \quad j, k, l = 1, \dots, N;$$

– центральные моменты четвертого порядка:

$$\mathbf{E}\{(\tilde{x}_j - \mu_j)(\tilde{x}_k - \mu_k)(\tilde{x}_l - \mu_l)(\tilde{x}_m - \mu_m)\} = \sigma_{jk}\sigma_{lm} + \sigma_{jl}\sigma_{km} + \sigma_{jm}\sigma_{lk}; \quad (4.4)$$

$$\mathbf{E}\{(\tilde{x}_j - \mu_j)^4\} = 3\sigma_{jj}^2, \quad j, k, l, m = 1, \dots, N.$$

– характеристическую функцию.

Теорема 4.1 (О характеристической функции многомерного нормального распределения). Пусть случайный N -мерный вектор $x \in \mathbb{R}^N$ имеет невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$), тогда его характеристическая функция

$$f_x(t) = \mathbf{E}\{e^{it^T x}\}, \quad t \in \mathbb{R}^N,$$

имеет вид

$$f_x(t) = e^{it^T \mu - \frac{1}{2}t^T \Sigma t}, \quad (4.5)$$

где i – мнимая единица.

Доказательство. Характеристическая функция непрерывного случайного N -мерного вектора x согласно формуле (2.11) определяется равенством:

$$f_x(t) = \int_{\mathbb{R}^N} \exp\left\{i \sum_{j=1}^N t_j z_j\right\} p_x(z) dz, \quad t, z \in \mathbb{R}^N, \quad (4.6)$$

где $p_x(z)$ – N -мерная плотность распределения вероятностей x .

Следовательно, если $\mathcal{L}\{x\} = \mathcal{N}_N(\mu, \Sigma)$, то, заменив в (4.20) $p_x(z)$ на $n_N(z|\mu, \Sigma)$ из (6.24), получаем:

$$\begin{aligned} f_x(t) &= \int_{\mathbb{R}^N} \exp\{it^T z\} n_N(z|\mu, \Sigma) dz = \\ &= \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \int_{\mathbb{R}^N} \exp\left\{it^T z - \frac{1}{2}(z - \mu)^T \Sigma^{-1}(z - \mu)\right\}, t, z \in \mathbb{R}^N. \end{aligned} \quad (4.7)$$

Выделим полный квадрат в выражении, стоящем под экспонентой:

$$\begin{aligned} it^T z - \frac{1}{2}(z - \mu)^T \Sigma^{-1}(z - \mu) &= \\ &= -\frac{1}{2}(z - (\mu + i\Sigma t) + i\Sigma t)^T \Sigma^{-1}(z - (\mu + i\Sigma t) + i\Sigma t) + it^T z = \\ &= -\frac{1}{2}(z - (\mu + i\Sigma t))^T \Sigma^{-1}(z - (\mu + i\Sigma t)) + it^T \mu - \frac{1}{2}t^T \Sigma t. \end{aligned} \quad (4.8)$$

Подставим (4.22) в (4.21):

$$\begin{aligned} f_x(t) &= \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \int_{\mathbb{R}^N} \exp\left\{-\frac{1}{2}(z - (\mu + i\Sigma t))^T \Sigma^{-1}(z - (\mu + i\Sigma t)) + it^T \mu - \frac{1}{2}t^T \Sigma t\right\} = \\ &= \exp\left\{it^T \mu - \frac{1}{2}t^T \Sigma t\right\} \int_{\mathbb{R}^N} \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(z - (\mu + i\Sigma t))^T \Sigma^{-1}(z - (\mu + i\Sigma t))\right\} = \\ &= \exp\left\{it^T \mu - \frac{1}{2}t^T \Sigma t\right\} \int_{\mathbb{R}^N} n_N(z|\mu + i\Sigma t, \Sigma) dz = \\ &= \exp\left\{it^T \mu - \frac{1}{2}t^T \Sigma t\right\}, t, z \in \mathbb{R}^N, \end{aligned} \quad (4.9)$$

так как в выполняется условие нормировки:

$$\int_{\mathbb{R}^N} n_N(z|\mu + i\Sigma t, \Sigma) dz = 1.$$

□

4.2 Линейные преобразования гауссовских случайных векторов

Теорема 4.2 (О линейном преобразовании гауссовского случайного вектора). Пусть $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$ – случайный вектор, имеющий невырожденное нормальное распределение:

$$\mathcal{L}\{x\} = \mathcal{N}_N(\mu, \Sigma),$$

C – неслучайная $(m \times N)$ -матрица

$$C = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1N} \\ c_{21} & c_{22} & \dots & c_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{m1} & c_{m2} & \dots & c_{mN} \end{pmatrix},$$

$m \leq N$, полного ранга ($\text{rank}(C) = m$), $a \in \mathbf{R}^m$ – фиксированный произвольный m -вектор. Тогда случайный вектор

$$y = Cx + a \quad (4.10)$$

является гауссовским случайным вектором, имеющим распределение $\mathcal{N}_N(\nu, \Xi)$, где

$$\begin{aligned}\nu &= C\mu + a, \\ \Xi &= C\Sigma C^T.\end{aligned}\tag{4.11}$$

Доказательство. Найдем характеристическую функцию $f_y(t)$ случайного вектора y , определенного линейным преобразованием (4.24) гауссовского случайного вектора x :

$$\begin{aligned}f_y(t) &= \mathbf{E} \left\{ e^{it^T y} \right\} = \mathbf{E} \left\{ e^{it^T (Cx+a)} \right\} = e^{it^T a} \mathbf{E} \left\{ e^{i(C^T t)^T x} \right\} = f_x(C^T t) = \\ &= e^{i(C^T t)^T \mu - \frac{1}{2}(C^T t)^T \Sigma C^T t} e^{it^T a} = e^{it^T (C\mu+a) - \frac{1}{2}t^T (C\Sigma C^T)t}.\end{aligned}\tag{4.12}$$

В правой части соотношения (4.25) получено выражение, которое представляет собой характеристическую функцию гауссовского m -мерного случайного вектора с распределением $\mathcal{N}_m(\nu, \Xi)$, где ν и Ξ определены в (4.23). Распределение вероятностей $\mathcal{N}_m(\nu, \Xi)$ невырождено, так как ковариационная $(m \times m)$ -матрица $\Xi = C\Sigma C^T$ имеет ранг m , в силу невырожденности исходной ковариационной $(N \times N)$ -матрицы Σ и полного ранга матрицы C ($\text{rank}(C) = m$) имеет ранг m и, следовательно, невырождена. \square

Теорема 4.3 (О маргинальных распределениях гауссовского случайного вектора). Пусть случайный N -вектор $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$ имеет невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$), и пусть $x^* = (\tilde{x}_{j_1}, \dots, \tilde{x}_{j_m})^T \in \mathbb{R}^m$ — m -вектор, образованный из каких-либо m ($1 \leq m \leq N$) компонент вектора x ($j_k \neq j_l \in \{1, \dots, N\}$, $k \neq l$, $k, l = 1, \dots, m$). Тогда случайный m -вектор x^* имеет невырожденное m -мерное маргинальное нормальное распределение с математическим ожиданием $\mu^* = (\tilde{\mu}_{j_k})_{k=1}^m \in \mathbb{R}^m$ и ковариационной $(m \times m)$ -матрицей $\Sigma^* = (\sigma_{j_k, j_l})_{k, l=1}^m$ ($|\Sigma^*| \neq 0$), образованными из соответствующих компонент исходных математического ожидания $\mu = (\tilde{\mu}_k)_{k=1}^N \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицы $\Sigma = (\sigma_{kl})_{k, l=1}^N$:

$$\mathcal{L}\{x^*\} = \mathcal{N}_m(\mu^*, \Sigma^*).$$

Доказательство. Представим m -вектор $x^* = (\tilde{x}_{j_1}, \dots, \tilde{x}_{j_m})^T$ как линейное преобразование исходного N -вектора $x \in \mathbb{R}^N$:

$$x^* = Cx,$$

где $(m \times N)$ -матрица C содержит единицы на пересечении k -й строки и j_k -го столбца ($j_k \in \{1, \dots, N\}$, $k = 1, \dots, m$), а остальные ее элементы равны нулю. Очевидно, что $\text{rank}(C) = m$, и из теоремы 4.2, с учетом вида матрицы C , получаем доказываемое. \square

Следствие 4.1. Компоненты случайного N -мерного вектора $x = (\tilde{x}_j)_{j=1}^N \in \mathbb{R}^N$, распределенного по невырожденному многомерному нормальному закону $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$), имеют одномерные маргинальные нормальные распределения:

$$\mathcal{L}\{\tilde{x}_j\} = \mathcal{N}_1(\tilde{\mu}_j, \sigma_{jj}), \quad \sigma_{jj} > 0, \quad j = 1, \dots, N.$$

Замечание 4.1.2. Если компоненты случайного вектора имеют одномерные нормальные распределения, то их совместное распределение, вообще говоря, не является многомерным нормальным. Оно будет таковым в случае их независимости (при этом некоррелированности компонент не достаточно!).

4.3 Независимость компонент гауссовских случайных векторов

Рассмотрим гауссовский случайный вектор $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$, распределенный по многомерному нормальному закону $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$). Предположим, что x разбит на K векторов-блоков:

$$\begin{aligned} x &= ((x^{(1)})^T : \dots : (x^{(K)})^T)^T = \\ &= \underbrace{(\tilde{x}_1, \dots, \tilde{x}_{N_1})^T}_{(x^{(1)})^T} \underbrace{(\tilde{x}_{N_1+1}, \dots, \tilde{x}_{N_1+N_2})^T}_{(x^{(2)})^T} \dots \underbrace{(\tilde{x}_{N_1+N_2+\dots+N_{K-1}+1}, \dots, \tilde{x}_N)^T}_{(x^{(K)})^T}, \end{aligned} \quad (4.13)$$

Очевидно, что размерность подвектора $x^{(k)}$ равна N_k , $k = 1, \dots, K$, и $N_1 + N_2 + \dots + N_K = N$.

Часто в прикладных задачах требуется установить зависимость или независимость векторов-блоков этого разбиения.

Теорема 4.4. Векторы-блоки $x^{(1)}, x^{(2)}, \dots, x^{(K)}$ разбиения (4.26) гауссовского случайного N -вектора x :

$$\mathcal{L}\{x\} = \mathcal{N}_N(\mu, \Sigma),$$

независимы тогда и только тогда, когда они некоррелированы.

Доказательство. Так как случайные векторы $x^{(1)}, x^{(2)}, \dots, x^{(K)}$ составлены из компонент гауссовского случайного вектора x , то согласно теореме 4.3, они также являются гауссовскими, причем:

$$\mathcal{L}\{x^{(k)}\} = \mathcal{N}_{N(k)}(\mu^{(k)}, \Sigma^{(k)}), \quad |\Sigma^{(k)}| \neq 0, \quad k = 1, \dots, K,$$

Но по известному критерию независимости [11] $x^{(1)}, x^{(2)}, \dots, x^{(K)}$ независимы тогда и только тогда, когда их совместная плотность распределения представима в виде произведения их плотностей, то есть

$$n_N(x|\mu, \Sigma) = \prod_{k=1}^K n_{N(k)}(x^{(k)}|\mu^{(k)}, \Sigma^{(k)}), \quad x = ((x^{(1)})^T : \dots : (x^{(K)})^T)^T \in \mathbb{R}^N.$$

Исходя из вида (6.24) плотности многомерного нормального распределения, заключаем, что последнее соотношение выполняется тогда и только тогда, когда случайные векторы $x^{(1)}, x^{(2)}, \dots, x^{(K)}$ некоррелированы: $\Sigma = \text{diag}\{\Sigma^{(1)}, \dots, \Sigma^{(K)}\}$. \square

Следствие 4.2. Пусть случайные векторы $x^{(1)}, x^{(2)}, \dots, x^{(K)}$ независимы в совокупности и имеют невырожденные нормальные распределения:

$$\mathcal{L}\{x^{(k)}\} = \mathcal{N}_{N(k)}(\mu^{(k)}, \Sigma^{(k)}), \quad |\Sigma^{(k)}| \neq 0, \quad k = 1, \dots, K.$$

Тогда составной вектор $x = ((x^{(1)})^T : \dots : (x^{(K)})^T)^T \in \mathbb{R}^N$, $N_1 + N_2 + \dots + N_K = N$, также имеет невырожденное нормальное распределение:

$$\mathcal{L}\{x\} = \mathcal{N}_N(\mu, \Sigma),$$

с вектором математического ожидания

$$\mu = ((\mu^{(1)})^T : \dots : (\mu^{(K)})^T)^T$$

и ковариационной матрицей

$$\Sigma = \text{diag}\{\Sigma^{(1)}, \dots, \Sigma^{(K)}\}, \quad |\Sigma| \neq 0.$$

4.4 Эллипсоид рассеяния

С невырожденным многомерным нормальным распределением $\mathcal{N}_N(\mu, \Sigma)$ в силу вида его плотности (6.24) связывают *метрику Махаланобиса*:

$$\rho(x, z) = \sqrt{(x - z)^T \Sigma^{-1} (x - z)}, \quad x, z \in \mathbb{R}^N,$$

которая позволяет построить в \mathbb{R}^N так называемый *эллипсоид рассеяния* с центром в точке $\mu \in \mathbb{R}^N$:

$$V_r = \{x : \rho(x, \mu) \leq r, \quad r > 0\}.$$

Воспользуемся известным в матричном анализе [9, 5] представлением для невырожденной ковариационной $(N \times N)$ -матрицы Σ , которая при этом в силу известного свойства неотрицательной определенности ковариационной матрицы положительно определена ($\Sigma \succ 0$):

$$\Sigma = \Sigma^{1/2} (\Sigma^{1/2})^T,$$

где невырожденная $(N \times N)$ -матрица $\Sigma^{1/2}$ определяется неоднозначно и является решением по Y матричного уравнения

$$Y^T \Sigma^{-1} Y = \mathbf{I}_N.$$

В частности, можно выбрать

$$\Sigma^{1/2} = (\sqrt{\lambda_1} \Psi_1 : \dots : \sqrt{\lambda_N} \Psi_N),$$

где $\{\lambda_k, \Psi_k\}_{k=1}^N$ – собственные числа и соответствующие им ортонормированные собственные векторы матрицы Σ , определенные в (3.24), (3.25).

Введем в рассмотрение случайный N -вектор

$$y = (\Sigma^{1/2})^{-1} (x - \mu).$$

Согласно теореме 4.1,

$$\mathcal{L}\{y\} = \mathcal{N}_N(\mathbf{0}_N, \mathbf{I}_N),$$

и случайная величина

$$\rho^2(x, \mu) = y^T y$$

является суммой N независимых в совокупности стандартных нормальных случайных величин и имеет χ^2 -распределение с N степенями свободы. Поэтому вероятность попадания случайного вектора x в эллипсоид рассеяния V_r радиуса r равна

$$\mathbf{P}\{x \in V_r\} = F_{\chi_N^2}(\mathbb{R}^2).$$

Этот факт позволяет по заданному уровню вероятности

$$\mathbf{P}\{x \in V_r\} = 1 - \alpha,$$

где $\alpha \in (0, 1)$ – мало, определить радиус эллипсоида рассеяния:

$$r = r(\alpha) = \sqrt{F_{\chi_N^2}^{-1}(1 - \alpha)},$$

где $F_{\chi_N^2}^{-1}(1 - \alpha)$ – квантиль уровня $1 - \alpha$ от χ^2 -распределения с N степенями свободы.

Эллипсоид рассеяния V_r , $r = r(\alpha)$, является областью концентрации наблюдений вокруг своего математического ожидания, а область в \mathbb{R}^N вне эллипсоида рассеяния содержит наблюдения, описываемые «хвостом» распределения вероятностей: вероятность попадания в эту область мала и равна α .

4.5 Условные распределения гауссовских случайных векторов. Частный коэффициент корреляции

Пусть случайный N -мерный вектор $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$ имеет невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$. Иногда необходимо исследовать зависимость определенного подмножества из m ($m < N$) компонент вектора x между собой при фиксированных (выбранных) значениях остальных $N - m$ компонент.

Предположим, не ограничивая общности, что исследуемыми m компонентами являются $\tilde{x}_1, \dots, \tilde{x}_m$. Это предположение порождает разбиение вектора $x \in \mathbb{R}^N$ его на два подвектора:

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix},$$

где

$$x^{(1)} = (\tilde{x}_1, \dots, \tilde{x}_m)^T \in \mathbb{R}^m, \quad x^{(2)} = (\tilde{x}_{m+1}, \dots, \tilde{x}_N)^T \in \mathbb{R}^{N-m},$$

а соответствующего ему N -вектора математического ожидания $\mu \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицы Σ – на соответствующие блоки:

$$\mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}, \quad \Sigma_{21} = \Sigma_{12}^T. \quad (4.14)$$

Согласно теореме 4.3 случайные векторы $x^{(1)}$ и $x^{(2)}$ являются гауссовскими с невырожденными нормальными распределениями:

$$\begin{aligned} \mathcal{L}\{x^{(1)}\} &= \mathcal{N}_m(\mu^{(1)}, \Sigma_{11}), \\ \mathcal{L}\{x^{(2)}\} &= \mathcal{N}_{N-m}(\mu^{(2)}, \Sigma_{22}). \end{aligned} \quad (4.15)$$

Вычислим условное распределение вероятностей случайного m -вектора $x^{(1)}$ при фиксированном значении компонент $(N - m)$ -вектора $x^{(2)}$.

Теорема 4.5 (Об условном распределении вероятностей гауссовских случайных векторов). Пусть случайный N -вектор $x \in \mathbb{R}^N$, имеющий невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$:

$$\mathcal{L}\{x\} = \mathcal{N}_N(\mu, \Sigma),$$

разбит на два подвектора $x^{(1)} \in \mathbb{R}^m$ и $x^{(2)} \in \mathbb{R}^{N-m}$:

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix}, \quad 1 \leq m < N.$$

Тогда условное распределение случайного вектора $x^{(1)}$ при фиксированном значении вектора $x^{(2)}$ – m -мерное нормальное со следующими значениями параметров:

$$\mathcal{L}\{x^{(1)}|x^{(2)}\} = \mathcal{N}_m(\mu_{1|2}, \Sigma_{11|2}), \quad (4.16)$$

$$\begin{aligned} \mu_{1|2} &= \mu^{(1)} + \Sigma_{12}\Sigma_{22}^{-1}(x^{(2)} - \mu^{(2)}), \\ \Sigma_{11|2} &= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T. \end{aligned} \quad (4.17)$$

Доказательство. Перейдем от случайного N -вектора-наблюдения

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix}, \quad x^{(1)} \in \mathbb{R}^m, \quad x^{(2)} \in \mathbb{R}^{N-m}, \quad 1 \leq m < N,$$

используя невырожденное линейное функциональное преобразование

$$y = y(x) : \mathbb{R}^N \rightarrow \mathbb{R}^N$$

следующего вида:

$$\begin{aligned} y^{(1)} &= x^{(1)} + Bx^{(2)}, \\ y^{(2)} &= x^{(2)}, \end{aligned} \tag{4.18}$$

к вспомогательному вектору

$$y = \begin{pmatrix} y^{(1)} \\ y^{(2)} \end{pmatrix}$$

где $(m \times (N - m))$ -матрица преобразования B выбирается так, чтобы случайные векторы $y^{(1)}$ и $y^{(2)}$ были некоррелированными, то есть, чтобы выполнялось равенство:

$$\mathbf{Cov}\{y^{(1)}, y^{(2)}\} = \mathbf{0}_{m \times (N-m)}, \tag{4.19}$$

где $\mathbf{0}_{m \times (N-m)}$ — нулевая $(m \times (N - m))$ -матрица.

Чтобы выбрать матрицу B , вычислим ковариацию между $y^{(1)}$ и $y^{(2)}$:

$$\mathbf{Cov}\{y^{(1)}, y^{(2)}\} = \mathbf{E}\left\{\left(y^{(1)} - \mathbf{E}\{y^{(1)}\}\right)\left(y^{(2)} - \mathbf{E}\{y^{(2)}\}\right)^{\mathbf{T}}\right\}, \tag{4.20}$$

Учитывая равенства (4.18) получаем:

$$\begin{aligned} \mathbf{E}\{y^{(1)}\} &= \mu^{(1)} + B\mu^{(2)}, \\ \mathbf{E}\{y^{(2)}\} &= \mu^{(2)}. \end{aligned}$$

Следовательно,

$$\begin{aligned} y^{(1)} - \mathbf{E}\{y^{(1)}\} &= x^{(1)} + Bx^{(2)} - \mu^{(1)} - B\mu^{(2)} = \dot{x}^{(1)} + B\dot{x}^{(2)}, \\ y^{(2)} - \mathbf{E}\{y^{(2)}\} &= x^{(2)} - \mu^{(2)} = \dot{x}^{(2)}, \end{aligned}$$

где

$$\begin{aligned} \dot{x}^{(1)} &::= x^{(1)} - \mu^{(1)}, \\ \dot{x}^{(2)} &::= x^{(2)} - \mu^{(2)}. \end{aligned}$$

Таким образом,

$$\begin{aligned} \mathbf{Cov}\{y^{(1)}, y^{(2)}\} &= \mathbf{E}\left\{\left(\dot{x}^{(1)} + B\dot{x}^{(2)}\right)\left(\dot{x}^{(2)}\right)^{\mathbf{T}}\right\} = \mathbf{E}\left\{\dot{x}^{(1)}\left(\dot{x}^{(2)}\right)^{\mathbf{T}}\right\} + B\mathbf{E}\left\{\dot{x}^{(2)}\left(\dot{x}^{(2)}\right)^{\mathbf{T}}\right\}, \\ &= \Sigma_{12} + B\Sigma_{22}. \end{aligned} \tag{4.21}$$

Приравнивая правую часть соотношения (4.21) к $\mathbf{0}_{m \times (N-m)}$, получаем

$$B = -\Sigma_{12}\Sigma_{22}^{-1}, \tag{4.22}$$

где в силу невырожденности исходной ковариационной матрицы Σ ее диагональный блок Σ_{22} также невырожден.

Так как случайный вектор

$$y = \begin{pmatrix} y^{(1)} \\ y^{(2)} \end{pmatrix}$$

является линейным преобразованием случайного вектора x , имеющего невырожденное нормальное распределение, то в силу теоремы 4.2, y также имеет нормальное распределение, а, согласно теореме 4.3, его подвекторы $y^{(1)}$ и $y^{(2)}$ – маргинальные нормальные распределения:

$$\begin{aligned}\mathcal{L}\{y^{(1)}\} &= \mathcal{N}_m(\mu^{(1)} - \Sigma_{12}\Sigma_{22}^{-1}\mu^{(2)}, \Sigma_{11|2}), \\ \mathcal{L}\{y^{(2)}\} &= \mathcal{N}_{N-m}(\mu^{(2)}, \Sigma_{22}),\end{aligned}\tag{4.23}$$

где $\Sigma_{11|2}$ – определена в (4.17) и является ковариационной матрицей вектора

$$y^{(1)} = x^{(1)} - \Sigma_{12}\Sigma_{22}^{-1}x^{(2)}.$$

Из некоррелированности случайных векторов $y^{(1)}$ и $y^{(2)}$ в силу их совместного нормального распределения следует их независимость. Поэтому плотность распределения вероятностей вектора y равна произведению плотностей распределения векторов $y^{(1)}$ и $y^{(2)}$:

$$p_y(u) = p_{y^{(1)}}(u^{(1)}) \cdot p_{y^{(2)}}(u^{(2)}) = n_m(u^{(1)}|\mu^{(1)} - \Sigma_{12}\Sigma_{22}^{-1}\mu^{(2)}, \Sigma_{11|2}) \cdot n_{N-m}(u^{(2)}|\mu^{(2)}, \Sigma_{22}),\tag{4.24}$$

где

$$u = \begin{pmatrix} u^{(1)} \\ u^{(2)} \end{pmatrix} \in \mathbf{R}^N, \quad u^{(1)} \in \mathbf{R}^m, \quad u^{(2)} \in \mathbf{R}^{N-m}.$$

Запишем теперь плотность распределения вектора x как преобразования y в x ($x = x(y)$), обратного преобразованию (4.18):

$$\begin{aligned}x^{(1)} &= y^{(1)} + \Sigma_{12}\Sigma_{22}^{-1}y^{(2)}, \\ x^{(2)} &= y^{(2)}.\end{aligned}$$

Отметим, что преобразования здесь имеют единичный якобиан, и из известного соотношения для плотности при невырожденном преобразовании [11] имеем:

$$\begin{aligned}p_x(v) &= p_y(u)|_{u=y(v)} = n_m(v^{(1)} - \Sigma_{12}\Sigma_{22}^{-1}v^{(2)}|\mu^{(1)} - \Sigma_{12}\Sigma_{22}^{-1}\mu^{(2)}, \Sigma_{11|2}) \cdot n_{N-m}(v^{(2)}|\mu^{(2)}, \Sigma_{22}) = \\ &= n_m(v^{(1)}|\mu_{1|2}, \Sigma_{11|2}) \cdot n_{N-m}(v^{(2)}|\mu^{(2)}, \Sigma_{22}),\end{aligned}\tag{4.25}$$

где

$$v = \begin{pmatrix} v^{(1)} \\ v^{(2)} \end{pmatrix} \in \mathbf{R}^N, \quad v^{(1)} \in \mathbf{R}^m, \quad v^{(2)} \in \mathbf{R}^{N-m}.$$

а m -вектор $\mu_{1|2}$ определен в (4.17).

С другой стороны, случайный вектор

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix}$$

имеет N -мерное распределение вероятностей с плотностью

$$p_x(v) = n_N(v|\mu, \Sigma),$$

а его подвектор $x^{(2)}$ – маргинальное нормальное распределение с плотностью

$$p_{x^{(2)}}(v^{(2)}) = n_{N-m}(v^{(2)}|\mu^{(2)}, \Sigma_{22}).$$

По известной формуле произведения плотностей [11] имеем

$$\begin{aligned}p_x(v) &= p_{x^{(1)}|x^{(2)}}(v^{(1)}|x^{(2)} = v^{(2)}) \cdot p_{x^{(2)}}(v^{(2)}) = \\ &= p_{x^{(1)}|x^{(2)}}(v^{(1)}|x^{(2)} = v^{(2)}) \cdot n_{N-m}(v^{(2)}|\mu^{(2)}, \Sigma_{22}).\end{aligned}\tag{4.26}$$

Приравнивая правые части соотношений (4.25) и (4.26), находим условную плотность:

$$p_{x^{(1)}|x^{(2)}}(v^{(1)}|x^{(2)} = v^{(2)}) = n_m(v^{(1)}|\mu_{1|2}, \Sigma_{11|2}).$$

Таким образом,

$$\mathcal{L}\{x^{(1)}|x^{(2)}\} = \mathcal{N}_m(\mu_{1|2}, \Sigma_{11|2}).$$

□

Определение 4.2. $(m \times m)$ -матрицу $\Sigma_{11|2}$, определенную в (4.17), иногда называют условной ковариационной матрицей, ее внедиагональные элементы – условными ковариациями, а диагональные элементы – условными дисперсиями.

В условиях теоремы 4.1 ковариационная $(m \times m)$ -матрица

$$\Sigma_{11|2} = \mathbf{Cov}\{x^{(1)}, x^{(1)}|x^{(2)}\}$$

из (4.17) может быть записана поэлементно:

$$\Sigma_{11|2} = \begin{pmatrix} \sigma_{11|m+1,\dots,N} & \sigma_{12|m+1,\dots,N} & \dots & \sigma_{1m|m+1,\dots,N} \\ \sigma_{21|m+1,\dots,N} & \sigma_{22|m+1,\dots,N} & \dots & \sigma_{2m|m+1,\dots,N} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1|m+1,\dots,N} & \sigma_{m2|m+1,\dots,N} & \dots & \sigma_{mm|m+1,\dots,N} \end{pmatrix},$$

где $\sigma_{kl|m+1,\dots,N}$ – условная ковариация k -й и l -й компонент при фиксированных значениях остальных компонент $x^{(2)} = (\tilde{x}_{m+1}, \dots, \tilde{x}_N)^T$:

$$\sigma_{kl|m+1,\dots,N} = \mathbf{Cov}\{\tilde{x}_k, \tilde{x}_l|x^{(2)}\} = \sigma_{kl} - \sigma_{(k)}^T \Sigma_{22}^{-1} \sigma_{(l)}, \quad k, l = 1, \dots, m, \quad (4.27)$$

а $\sigma_{(k)}^T$ – k -я строка матрицы Σ_{12} :

$$\sigma_{(k)}^T = \mathbf{Cov}\{\tilde{x}_k, x_{(2)}\}, \quad k = 1, \dots, m. \quad (4.28)$$

Определение 4.3. Частным коэффициентом корреляции между k -й и l -й компонентами \tilde{x}_k и \tilde{x}_l ($k, l = 1, \dots, m$) случайного N -вектора x при фиксированных значениях компонент

$$x^{(2)} = (\tilde{x}_{m+1}, \dots, \tilde{x}_N)^T$$

называется величина

$$\rho_{kl|m+1,\dots,N} = \frac{\sigma_{kl|m+1,\dots,N}}{\sqrt{\sigma_{kk|m+1,\dots,N} \sigma_{ll|m+1,\dots,N}}}, \quad (4.29)$$

где условная ковариация $\sigma_{kl|m+1,\dots,N}$ называется частной ковариацией, а условная дисперсия $\sigma_{kk|m+1,\dots,N}$ – частной дисперсией.

Частный коэффициент корреляции является характеристикой парной взаимосвязи на множестве выбранных компонент при фиксированных значениях остальных компонент и не зависит от фиксируемых значений компонент.

Отметим, что частные корреляция, ковариация и дисперсия обладают всеми свойствами своих обычных (безусловных) аналогов.

4.6 Функция регрессии. Множественный коэффициент корреляции

Пусть случайный N -мерный вектор $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T$ разбит на два подвектора:

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix},$$

где

$$x^{(1)} = (\tilde{x}_1, \dots, \tilde{x}_m)^T \in \mathbb{R}^m, \quad x^{(2)} = (\tilde{x}_{m+1}, \dots, \tilde{x}_N)^T \in \mathbb{R}^{N-m}.$$

Определение 4.4. Функцией регрессии подвектора $x^{(1)} \in \mathbb{R}^m$ на $x^{(2)} \in \mathbb{R}^{N-m}$ называется условное математическое ожидание

$$\mu_{1|2} = \mu_{1|2}(x^{(2)}) = \mathbf{E} \{x^{(1)} | x^{(2)}\}. \quad (4.30)$$

Предположим, что N -вектор $x \in \mathbb{R}^N$ имеет невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$. В этом случае функция регрессии $x^{(1)}$ на $x^{(2)}$, согласно (4.17), линейна и имеет вид:

$$\mu_{1|2} = \mu^{(1)} + \Sigma_{12} \cdot \Sigma_{22}^{-1}(x^{(2)} - \mu^{(2)}),$$

где $\mu^{(1)}$, $\mu^{(2)}$, Σ_{12} и Σ_{22} определены в (6.24). Матрица $B = \Sigma_{12} \cdot \Sigma_{22}^{-1}$ называется *матрицей коэффициентов регрессии*:

$$B = \Sigma_{12} \cdot \Sigma_{22}^{-1} = \begin{pmatrix} \beta_{1,1} & \beta_{1,2} & \dots & \beta_{1,N-m} \\ \beta_{2,1} & \beta_{2,2} & \dots & \beta_{2,N-m} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{m,1} & \beta_{m,2} & \dots & \beta_{m,N-m} \end{pmatrix} = \begin{pmatrix} \beta_1^T \\ \beta_2^T \\ \vdots \\ \beta_m^T \end{pmatrix}, \quad (4.31)$$

где β_k^T – k -я строка матрицы $B = \Sigma_{12} \cdot \Sigma_{22}^{-1}$:

$$\beta_k^T = \sigma_{(k)}^T \cdot \Sigma_{22}^{-1} \quad k = 1, \dots, m, \quad (4.32)$$

а $\sigma_{(k)}^T$ – k -я строка матрицы Σ_{12} :

$$\sigma_{(k)}^T = (\sigma_{k,m+1} \sigma_{k,m+2} \dots \sigma_{k,N}) \quad k = 1, \dots, m. \quad (4.33)$$

Учитывая (4.32) и (4.33) функцию регрессии k й компоненты \tilde{x}_k вектора $x^{(1)}$ на $x^{(2)}$ можно записать в виде:

$$\begin{aligned} \mathbf{E} \{ \tilde{x}_k | x^{(2)} \} &= \mu_k + \beta_k^T \cdot (x^{(2)} - \mu^{(2)}) = \mu_k - \beta_k^T \cdot (\mu^{(2)} - x^{(2)}) = \\ &= \mu_k - \sum_{j=m+1}^N \beta_{k,j} (\mu_j - \tilde{x}_j) \quad k = 1, \dots, m. \end{aligned} \quad (4.34)$$

Замечание 4.4.1. Из формулы (4.34) видно, что коэффициент регрессии $\beta_{k,j}$ можно рассматривать как коэффициент влияния компоненты \tilde{x}_j на $\mathbf{E} \{ \tilde{x}_k | x^{(2)} \}$:

- чем больше значение $|\beta_{k,j}|$, тем сильнее влияет \tilde{x}_j на $\mathbf{E} \{ \tilde{x}_k | x^{(2)} \}$, $k = 1, \dots, m$, $j = m+1, \dots, N$;
- если $\beta_{k,j} = 0$, то \tilde{x}_j не влияет на $\mathbf{E} \{ \tilde{x}_k | x^{(2)} \}$, $k = 1, \dots, m$, $j = m+1, \dots, N$.

Предположим, что необходимо оценить неизвестные значения компонент из $x^{(1)} \in \mathbb{R}^m$ по наблюдаемым компонентам $x^{(2)} \in \mathbb{R}^{N-m}$ ($m < N$). Эта задача имеет большое прикладное значение при прогнозировании и оценивании («заполнении») пропущенных значений

компонент в выборке. Математически задача прогнозирования (или оценивания пропущенных значений компонент в выборке) заключается в построении функционального преобразования $g(\cdot) : \mathbb{R}^{N-m} \rightarrow \mathbb{R}^m$

$$\hat{x}^{(1)} = g(x^{(2)}), \quad g(\cdot) \in \mathcal{G} \subset \mathbb{R}^m, \quad (4.35)$$

где $\hat{x}^{(1)}$ – вычисленное предсказание (оценка) для $x^{(1)}$, называемое *прогнозом* или *значением прогноза*, а $g(\cdot)$ – *функция прогнозирования (предиктор)*, \mathcal{G} – некоторый допустимый класс функций.

Теорема 4.6 (Об оптимальных свойствах функции регрессии). *В условиях теоремы 4.5 среди всех прогнозов*

$$\hat{x}^{(1)} = (\hat{x}_1, \dots, \hat{x}_m)^T$$

для $x^{(1)}$ по $x^{(2)}$ прогноз

$$\hat{x}_*^{(1)} = (\hat{x}_1^*, \dots, \hat{x}_m^*)^T = \mu_{1|2}(x^{(2)}),$$

где $\mu_{1|2} = \mu_{1|2}(x^{(2)})$ – функция регрессии, определенная в (4.30), является оптимальным прогнозом в смысле среднеквадратической ошибки прогнозирования, то есть:

$$\hat{x}_k^* = \arg \min_{\hat{x}_k} \mathbf{E}\{(\hat{x}_k - \tilde{x}_k)^2\}, \quad k = 1, \dots, m, \quad (4.36)$$

и максимально коррелирован с предсказываемым значением \tilde{x}_k

$$\hat{x}_k^* = \arg \max_{\hat{x}_k} \mathbf{Corr}\{\hat{x}_k, \tilde{x}_k\}, \quad k = 1, \dots, m. \quad (4.37)$$

Доказательство. Пусть \hat{x}_k – любой допустимый прогноз компоненты \tilde{x}_k подвектора $x^{(1)}$. Запишем для этого прогноза среднеквадратическую ошибку прогнозирования $\Delta^2(\hat{x}_k)$:

$$\begin{aligned} \Delta^2(\hat{x}_k) &= \mathbf{E}\{(\hat{x}_k - \tilde{x}_k)^2\} = \mathbf{E}\{(\hat{x}_k - \hat{x}_k^* + \hat{x}_k^* - \tilde{x}_k)^2\} = \\ &= \mathbf{E}\{(\hat{x}_k - \hat{x}_k^*)^2\} + 2\mathbf{E}\{(\hat{x}_k - \hat{x}_k^*)(\hat{x}_k^* - \tilde{x}_k)\} + \mathbf{E}\{(\hat{x}_k^* - \tilde{x}_k)^2\}, \quad k = 1, \dots, m, \end{aligned} \quad (4.38)$$

где \hat{x}_k^* – условное математическое ожидание \tilde{x}_k при фиксированном значении вектора $x^{(1)}$:

$$\hat{x}_k^* = \mathbf{E}\{\tilde{x}_k | x^{(2)}\} \quad k = 1, \dots, m.$$

По формуле полного математического ожидания [11] имеем для любого допустимого прогноза \hat{x}_k :

$$\begin{aligned} \mathbf{E}\{(\hat{x}_k - \hat{x}_k^*)(\hat{x}_k^* - \tilde{x}_k)\} &= \mathbf{E}\left\{\mathbf{E}\{(\hat{x}_k - \hat{x}_k^*)(\hat{x}_k^* - \tilde{x}_k) | x^{(2)}\}\right\} = \\ &= \mathbf{E}\left\{(\hat{x}_k - \hat{x}_k^*)\mathbf{E}\{(\hat{x}_k^* - \tilde{x}_k) | x^{(2)}\}\right\} = \mathbf{E}\left\{(\hat{x}_k - \hat{x}_k^*)(\hat{x}_k^* - \hat{x}_k^*)\right\} = 0 \quad k = 1, \dots, m. \end{aligned}$$

Следовательно, так как

$$\mathbf{E}\{(\hat{x}_k^* - \tilde{x}_k)^2\} = \Delta^2(\hat{x}_k^*)$$

формулу (4.40) можно записать в виде:

$$\Delta^2(\hat{x}_k) = \mathbf{E}\{(\hat{x}_k - \hat{x}_k^*)^2\} + \Delta^2(\hat{x}_k^*) \geq \Delta^2(\hat{x}_k^*), \quad k = 1, \dots, m. \quad (4.39)$$

Равенство будет в том случае, когда $\hat{x}_k \equiv \hat{x}_k^*$, $k = 1, \dots, m$. Таким образом, \hat{x}_k^* – оптимальный прогноз в смысле среднеквадратической ошибки прогнозирования компоненты \tilde{x}_k подвектора $x^{(1)}$ при фиксированном $x^{(2)}$, $k = 1, \dots, m$.

Для доказательства второго утверждения теоремы вычислим ковариацию произвольного прогноза $\hat{\tilde{x}}_k$ с прогнозируемой компонентой \tilde{x}_k , учитывая, что

$$\mathbf{E}\{\tilde{x}_k\} = \mathbf{E}\{\mathbf{E}\{\tilde{x}_k|x^{(2)}\}\} = \mathbf{E}\{\hat{\tilde{x}}_k^*\}, \quad k = 1, \dots, m:$$

$$\begin{aligned} \mathbf{Cov}\{\hat{\tilde{x}}_k, \tilde{x}_k\} &= \mathbf{E}\{(\hat{\tilde{x}}_k - \mathbf{E}\{\hat{\tilde{x}}_k\})(\tilde{x}_k - \mathbf{E}\{\tilde{x}_k\})\} = \mathbf{E}\{(\hat{\tilde{x}}_k - \mathbf{E}\{\hat{\tilde{x}}_k\})\mathbf{E}\{(\tilde{x}_k - \mathbf{E}\{\tilde{x}_k\})|x^{(2)}\}\} = \\ &= \mathbf{E}\{(\hat{\tilde{x}}_k - \mathbf{E}\{\hat{\tilde{x}}_k\})(\hat{\tilde{x}}_k^* - \mathbf{E}\{\hat{\tilde{x}}_k^*\})\} = \mathbf{Cov}\{\hat{\tilde{x}}_k, \hat{\tilde{x}}_k^*\}, \quad k = 1, \dots, m. \end{aligned}$$

Следовательно,

$$\begin{aligned} \mathbf{Corr}^2\{\hat{\tilde{x}}_k, \tilde{x}_k\} &= \frac{\mathbf{Cov}^2\{\hat{\tilde{x}}_k, \tilde{x}_k\}}{\mathbf{D}\{\hat{\tilde{x}}_k\}\mathbf{D}\{\tilde{x}_k\}} = \frac{\mathbf{Cov}^2\{\hat{\tilde{x}}_k, \hat{\tilde{x}}_k^*\}}{\mathbf{D}\{\hat{\tilde{x}}_k\}\mathbf{D}\{\hat{\tilde{x}}_k^*\}} \frac{\mathbf{D}\{\hat{\tilde{x}}_k^*\}}{\mathbf{D}\{\tilde{x}_k\}} = \\ &= \mathbf{Corr}^2\{\hat{\tilde{x}}_k, \hat{\tilde{x}}_k^*\}\mathbf{Corr}^2\{\hat{\tilde{x}}_k^*, \tilde{x}_k\}, \quad k = 1, \dots, m. \end{aligned}$$

Поскольку $\mathbf{Corr}^2\{\hat{\tilde{x}}_k, \hat{\tilde{x}}_k^*\} \leq 1$, то получаем

$$|\mathbf{Corr}\{\hat{\tilde{x}}_k, \tilde{x}_k\}| \leq |\mathbf{Corr}\{\hat{\tilde{x}}_k^*, \tilde{x}_k\}|, \quad k = 1, \dots, m.$$

□

Следствие 4.3.

$$\mathbf{E}\{(\hat{\tilde{x}}_k^* - \tilde{x}_k)^2\} = \sigma_{kk|m+1, \dots, N}; \quad (4.40)$$

$$\mathbf{Corr}\{\hat{\tilde{x}}_k^*, \tilde{x}_k\} = \sqrt{\frac{\sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} \sigma_{(k)}}{\sigma_{kk}}}, \quad (4.41)$$

где Σ_{22} – ковариационная матрица вектора $x^{(2)}$, σ_{kk} – дисперсия k -й компоненты, а $(N - m)$ -вектор $\sigma_{(k)}$ – определен в (4.27), $k = 1, \dots, m$.

Доказательство. Для доказательства формул (4.40) и (4.41) воспользуемся соотношением (4.30) и запишем оптимальный прогноз $\hat{x}_*^{(1)} = (\hat{x}_1^*, \dots, \hat{x}_m^*)^{\mathbf{T}}$ покомпонентно:

$$\hat{x}_k^* = \tilde{\mu}_k + \sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)}), \quad k = 1, \dots, m, \quad (4.42)$$

где $\tilde{\mu}_k = \mathbf{E}\{\tilde{x}_k\}$, а $\mu^{(2)} = \mathbf{E}\{x^{(2)}\}$. Учтем, что

$$\mathbf{E}\{\hat{\tilde{x}}_k^*\} = \mathbf{E}\{\mathbf{E}\{\tilde{x}_k|x^{(2)}\}\} = \mathbf{E}\{\tilde{x}_k\} = \tilde{\mu}_k,$$

и вычислим ковариацию:

$$\begin{aligned} \mathbf{Cov}\{\hat{\tilde{x}}_k^*, \tilde{x}_k\} &= \mathbf{E}\{(\hat{\tilde{x}}_k^* - \tilde{\mu}_k)(\tilde{x}_k - \tilde{\mu}_k)\} = \\ &= \mathbf{E}\{\sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)})(\tilde{x}_k - \tilde{\mu}_k)\} = \sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} \sigma_{(k)}. \end{aligned}$$

Отметим, что $\mathbf{D}\{\tilde{x}_k\} = \sigma_{kk}$, а для $\mathbf{D}\{\hat{\tilde{x}}_k^*\}$ имеем

$$\begin{aligned} \mathbf{D}\{\hat{\tilde{x}}_k^*\} &= \mathbf{E}\{(\hat{\tilde{x}}_k^* - \tilde{\mu}_k)^2\} = \mathbf{E}\{(\sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} (x^{(2)} - \mu^{(2)}))^2\} = \\ &= \sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} \mathbf{Cov}\{x^{(2)}, x^{(2)}\} \Sigma_{22}^{-1} \sigma_{(k)} = \sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} \sigma_{(k)}. \end{aligned}$$

Полученные соотношения и позволяют вычислить корреляцию в (4.41), а также доказать (4.40):

$$\begin{aligned} \mathbf{E}\{(\hat{\tilde{x}}_k^* - \tilde{x}_k)^2\} &= \mathbf{D}\{\hat{\tilde{x}}_k^* - \tilde{x}_k\} = \mathbf{D}\{\hat{\tilde{x}}_k^*\} - 2\mathbf{Cov}\{\hat{\tilde{x}}_k^*, \tilde{x}_k\} + \mathbf{D}\{\tilde{x}_k\} = \\ &= \sigma_{kk} - \sigma_{(k)}^{\mathbf{T}} \Sigma_{22}^{-1} \sigma_{(k)} = \sigma_{kk|m+1, \dots, N}, \end{aligned}$$

где учтено (4.27).

□

Определение 4.5. Множественным коэффициентом корреляции *между компонентой* \tilde{x}_k ($k = 1, \dots, m$) *и компонентами подвектора* $x^{(2)} = (\tilde{x}_{m+1}, \dots, \tilde{x}_N)^T$ *называется величина*

$$\bar{R}_{k,m+1,\dots,N} = \sqrt{\frac{\sigma_{(k)}^T \Sigma_{22}^{-1} \sigma_{(k)}}{\sigma_{kk}}}. \quad (4.43)$$

Свойства множественного коэффициента корреляции

Свойство 1. $0 \leq \bar{R}_{k,m+1,\dots,N} \leq 1$, $k = 1, \dots, m$.

Доказательство. Доказательство следует из определения (4.43) и соотношения (4.27):

$$\mathbf{D}\{\tilde{x}_k | x^{(2)}\} = \sigma_{kk|m+1,\dots,N} = \sigma_{kk} - \sigma_{(k)}^T \Sigma_{22}^{-1} \sigma_{(k)} \geq 0, \quad k = 1, \dots, m.$$

□

Свойство 2. $\bar{R}_{k,m+1,\dots,N} = 0$ тогда и только тогда, когда k -я компонента \tilde{x}_k ($k = 1, \dots, m$) некоррелирована с компонентами из множества $\{\tilde{x}_{m+1}, \dots, \tilde{x}_N\}$.

Доказательство. Из (4.43) видно, что $\bar{R}_{k,m+1,\dots,N} = 0$ тогда и только тогда, когда

$$\sigma_{(k)} = \mathbf{Cov}\{x^{(2)}, \tilde{x}_k\} = \mathbf{0}_{N-m}, \quad k = 1, \dots, m.$$

□

Свойство 3. $\bar{R}_{k,m+1,\dots,N} = 1$ тогда и только тогда, когда k -я компонента \tilde{x}_k представима в виде

$$\tilde{x}_k \stackrel{\text{P} \equiv 1}{=} b_k^T x^{(2)} + d_k,$$

где $b_k \in \mathbb{R}^{N-m}$, $d_k \in \mathbb{R}^1$ – детерминированы, $k = 1, \dots, m$.

Доказательство. Из (4.43), (4.41) и свойств обычного коэффициента корреляции следует, что

$$\bar{R}_{k,m+1,\dots,N} = \mathbf{Corr}\{\hat{x}_k^*, \tilde{x}_k\} = 1$$

тогда и только тогда, когда

$$\tilde{x}_k \stackrel{\text{P} \equiv 1}{=} \tilde{b}_k^T \hat{x}_k^* + \tilde{d}_k,$$

где $\tilde{b}_k \in \mathbb{R}^{N-m}$, $\tilde{d}_k \in \mathbb{R}^1$ – детерминированы, $k = 1, \dots, m$. Подставим вместо \hat{x}_k^* его выражение через $x^{(2)}$ из (4.42) и получим доказываемое. □

Свойство 4. Справедливо следующее соотношение, связывающее частную дисперсию $\sigma_{kk|m+1,\dots,N}$ из (4.27) и множественный коэффициент корреляции $\bar{R}_{k,m+1,\dots,N}$ ($k = 1, \dots, m$):

$$\sigma_{kk|m+1,\dots,N} = (1 - \bar{R}_{k,m+1,\dots,N}^2) \sigma_{kk},$$

и частная дисперсия $\sigma_{kk|m+1,\dots,N}$ никогда не превосходит соответствующей безусловной дисперсии σ_{kk} : $\sigma_{kk|m+1,\dots,N} \leq \sigma_{kk}$.

Доказательство. непосредственно следует из соотношения (4.27) и свойства 1. □

Свойство 5. Если $N = 2$, $m = 1$, то множественный коэффициент корреляции (4.43) совпадает с точностью до знака с парным коэффициентом корреляции:

$$\bar{R}_{1,2} = |\rho_{12}|.$$

Замечание 4.5.1. В теории множественной регрессии множественный коэффициент корреляции служит мерой оптимальности прогноза: чем он ближе к единице, тем точнее прогноз.

5 Статистическое оценивание параметров многомерной гауссовской модели

5.1 Оценки максимального правдоподобия параметров многомерного нормального распределения

Пусть имеется случайная выборка $X = \{x_1, \dots, x_n\}$ объема n из невырожденного N -мерного нормального распределения $\mathcal{N}_N(\mu, \Sigma)$ с неизвестными истинными значениями параметров: математического ожидания $\mu \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицы Σ ($|\Sigma| \neq 0$), для построения статистических оценок которых ($\hat{\mu}$ и $\hat{\Sigma}$) по выборке X воспользуемся методом максимального правдоподобия [12].

Теорема 5.1 (Об оценках максимального правдоподобия параметров многомерного нормального распределения). *Пусть наблюдения $x_1, \dots, x_n \in \mathbb{R}^N$, образующие выборку X объема $n > N$, являются независимыми в совокупности, одинаково распределенными случайными N -векторами с невырожденным нормальным распределением $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$). Тогда единственными оценками максимального правдоподобия (МП-оценками) для вектора математического ожидания μ и ковариационной матрицы Σ являются соответственно выборочное среднее $\hat{\mu}$:*

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \in \mathbb{R}^N, \quad (5.1)$$

и выборочная ковариационная $(N \times N)$ -матрица $\hat{\Sigma}$:

$$\hat{\Sigma} = \frac{1}{n} A, \quad (5.2)$$

где

$$A = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T. \quad (5.3)$$

Доказательство. Плотность распределения вероятностей $n_N(\cdot|\mu, \Sigma)$ многомерного нормального закона распределения $\mathcal{N}_N(\mu, \Sigma)$, согласно определению (6.24) имеет вид:

$$n_N(z|\mu, \Sigma) = (2\pi)^{-\frac{N}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(z - \mu)^T \Sigma^{-1}(z - \mu)\right).$$

Следовательно, логарифмическая функция правдоподобия $l(\mu, \Sigma)$ [12] для неизвестных значений параметров μ и Σ , построенная по случайной выборке по выборке $X = \{x_1, \dots, x_n\}$, записывается следующим образом:

$$\begin{aligned} l(\mu, \Sigma) &= \ln \left(\prod_{i=1}^n n_N(x_i|\mu, \Sigma) \right) = \ln \left(\prod_{i=1}^n (2\pi)^{-\frac{N}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu)\right) \right) = \\ &= \sum_{i=1}^n \left(-\frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma| - \frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu) \right) = \\ &= -\frac{nN}{2} \ln(2\pi) - \frac{n}{2} \ln |\Sigma| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1}(x_i - \mu) = \\ &= -\frac{nN}{2} \ln(2\pi) + \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n \text{tr}(\Sigma^{-1}(x_i - \mu)(x_i - \mu)^T) = \end{aligned}$$

$$= -\frac{nN}{2} \ln(2\pi) + \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \text{tr}(\Sigma^{-1} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^{\mathbf{T}}), \quad (5.4)$$

где использованы следующие свойства следа матрицы $\text{tr}(\cdot)$:

- $c = \text{tr}(c)$, $c \in \mathbb{R}$;
- $\text{tr}(BC) = \text{tr}(CB)$;
- $\text{tr}(B + D) = \text{tr}(B) + \text{tr}(D)$.

В формуле (5.4) преобразуем сумму

$$\sum_{i=1}^n (x_i - \mu)(x_i - \mu)^{\mathbf{T}},$$

учитывая очевидное тождество

$$\sum_{t=1}^n (x_t - \bar{x}) \equiv \mathbf{0}_N. \quad (5.5)$$

Имеем:

$$\begin{aligned} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^{\mathbf{T}} &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)(x_i - \bar{x} + \bar{x} - \mu)^{\mathbf{T}} = \\ &= \sum_{i=1}^n [(x_i - \bar{x})(x_i - \bar{x})^{\mathbf{T}} + (x_i - \bar{x})(\bar{x} - \mu)^{\mathbf{T}} + (\bar{x} - \mu)(x_i - \bar{x})^{\mathbf{T}} + (\bar{x} - \mu)(\bar{x} - \mu)^{\mathbf{T}}] = \\ &= A + n(\bar{x} - \mu)(\bar{x} - \mu)^{\mathbf{T}}, \end{aligned} \quad (5.6)$$

где \bar{x} и A определены соответственно в (5.1) и (5.3).

Подставив (5.6) в (5.4) и используя свойства следа матрицы, получаем:

$$l(\mu, \Sigma) = -\frac{nN}{2} \ln(2\pi) + \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \text{tr}(\Sigma^{-1} A) - \frac{n}{2} (\bar{x} - \mu)^{\mathbf{T}} \Sigma^{-1} (\bar{x} - \mu). \quad (5.7)$$

Из вида (5.7) логарифмической функции правдоподобия $l(\mu, \Sigma)$ заключаем, что максимум по $\mu \in \mathbb{R}^N$ в ней достигается лишь на выборочном среднем $\hat{\mu} = \bar{x}$, которое и является МП-оценкой для вектора математического ожидания μ .

Из (5.7) также видно, что вместо МП-оценки $\hat{\Sigma}$ для ковариационной матрицы Σ целесообразно строить МП-оценку

$$\widehat{\Sigma^{-1}} = (\sigma_{kj}^*)_{k,j=1}^N$$

для обратной ковариационной матрицы Σ^{-1} . В силу взаимно-однозначного соответствия между Σ и Σ^{-1} МП-оценкой для Σ будет

$$\hat{\Sigma} = (\widehat{\Sigma^{-1}})^{-1},$$

которая находится из решения задачи:

$$l^*(\Sigma^{-1}) = l(\bar{x}, \Sigma) = -\frac{nN}{2} \ln(2\pi) + \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \text{tr}(\Sigma^{-1} A) \rightarrow \max_{\Sigma^{-1}}.$$

Вычислим частные производные от $l^*(\Sigma^{-1})$ по σ_{kj}^* , используя поэлементное представление матрицы A :

$$A = (a_{kj})_{k,j=1}^N,$$

разложение определителя по строке и известное в матричном анализе представление Крамера для элементов обратной матрицы [9, 5]:

$$\frac{\partial l^*(\Sigma^{-1})}{\partial \sigma_{kj}^*} = \frac{n}{2} \frac{1}{|\Sigma^{-1}|} \frac{\partial |\Sigma^{-1}|}{\partial \sigma_{kj}^*} - \frac{1}{2} \frac{\partial}{\partial \sigma_{kj}^*} \sum_{i,l=1}^N \sigma_{il}^* a_{li} = \frac{n}{2} \sigma_{kj} - \frac{1}{2} a_{jk}, \quad k, j = 1, \dots, N.$$

Согласно необходимому условию максимума, полученные выражения для производных приравняем к нулю и получаем МП-оценку $\hat{\Sigma} = (\hat{\sigma}_{kj})_{k,j=1}^N$ для ковариационной матрицы $\Sigma = (\sigma_{kj})_{k,j=1}^N$:

$$\hat{\sigma}_{kj} = \frac{1}{n} a_{jk}, \quad k, j = 1, \dots, N,$$

откуда в силу симметричности матрицы A :

$$a_{jk} = a_{kj}, \quad k, j = 1, \dots, N,$$

окончательно имеем

$$\hat{\Sigma} = \frac{1}{n} A,$$

что совпадает с (5.2). Условие $n > N$ обеспечивает невырожденность полученной оценки ковариационной матрицы: $\mathbf{P}\{|A| = 0\} = 0$. \square

Выборочное среднее

$$\hat{\mu} = \bar{x}$$

и выборочная ковариационная матрица

$$\hat{\Sigma} = \frac{1}{n} A,$$

полученные в теореме 5.1, как МП-оценки обладают всеми их свойствами: сильная состоятельность, асимптотическая несмещенность и асимптотическая нормальность.

5.2 Распределение вероятностей оценок максимального правдоподобия параметров многомерного нормального распределения

Пусть C – ортогональная $(n \times n)$ -матрица:

$$C = (c_{ij})_{i,j=1}^n, \quad (5.8)$$

обладающую свойством [9], [5]:

$$\begin{aligned} CC^T &= C^T C = \mathbf{I}_n, \quad (C^{-1} = C^T), \\ \sum_{m=1}^n c_{km} \cdot c_{ml} &= \delta_{kl} = \begin{cases} 1, & k = l, \\ 0, & k \neq l, \end{cases} \quad k, l = 1, \dots, n. \end{aligned} \quad (5.9)$$

Лемма 5.1. Пусть x_1, \dots, x_n – независимые в совокупности N -мерные гауссовские случайные векторы:

$$\mathcal{L}\{x_i\} = \mathcal{N}_N(\mu_i, \Sigma), \quad i = 1, \dots, n.$$

Тогда случайные векторы y_1, \dots, y_n :

$$y_i = \sum_{j=1}^n c_{ij} \cdot x_j \quad i = 1, \dots, n, \quad (5.10)$$

где $C = (c_{ij})_{i,j=1}^n$ – ортогональная матрица, определенная в (5.8) и (5.9), также являются независимыми и имеют нормальные распределения:

$$\mathcal{L}\{y_i\} = \mathcal{N}_N(\nu_i, \Sigma), \quad i = 1, \dots, n, \quad (5.11)$$

причем

$$\nu_i = \sum_{j=1}^n c_{ij} \cdot \mu_j \quad i = 1, \dots, n. \quad (5.12)$$

Кроме того, остается неизменной матричная сумма квадратов:

$$\sum_{j=1}^n x_j \cdot (x_j)^{\mathbf{T}} = \sum_{j=1}^n y_j \cdot (y_j)^{\mathbf{T}}. \quad (5.13)$$

Доказательство. Случайные векторы y_1, \dots, y_n являются линейными преобразованиями нормально распределенных независимых в совокупности случайных векторов x_1, \dots, x_n , имеющих, согласно следствию 4.2, совместное нормальное распределение. По теореме 4.2 y_1, \dots, y_n также имеют совместное нормальное распределение и по теореме 4.3 нормально распределены:

$$\mathcal{L}\{y_i\} = \mathcal{N}_N(\nu_i, \Sigma_i), \quad i = 1, \dots, n.$$

Найдем параметры ν_i и Σ_i , $i = 1, \dots, n$, этих распределений.

Для математических ожиданий имеем

$$\nu_i = \mathbf{E}\{y_i\} = \mathbf{E}\left\{\sum_{j=1}^n c_{ij}x_j\right\} = \sum_{j=1}^n c_{ij}E\{x_j\} = \sum_{j=1}^n c_{ij}\mu_j, \quad i = 1, \dots, n.$$

Вместо вычисления ковариационных матриц Σ_i , $i = 1, \dots, n$, решим более общую задачу – найдем ковариации векторов y_1, \dots, y_n :

$$\begin{aligned} \mathbf{Cov}\{y_i, y_l\} &= \mathbf{E}\{(y_i - \nu_i)(y_l - \nu_l)^{\mathbf{T}}\} = \\ &= \mathbf{E}\left\{\left(\sum_{j=1}^n c_{ij}x_j - \sum_{j=1}^n c_{ij}\mu_j\right)\left(\sum_{k=1}^n c_{lk}x_k - \sum_{k=1}^n c_{lk}\mu_k\right)^{\mathbf{T}}\right\} = \\ &= \sum_{j=1}^n \sum_{k=1}^n c_{ij}c_{lk}\mathbf{E}\{(x_j - \mu_j)(x_k - \mu_k)^{\mathbf{T}}\} = \sum_{j=1}^n \sum_{k=1}^n c_{ij}c_{lk}\mathbf{Cov}\{x_j, x_k\}. \end{aligned} \quad (5.14)$$

Так как случайные векторы x_1, \dots, x_n – независимы в совокупности и

$$\mathbf{Cov}\{x_j, x_j\} = \Sigma, \quad j = 1, \dots, n,$$

то

$$\mathbf{Cov}\{x_j, x_k\} = \delta_{jk}\Sigma = \begin{cases} \Sigma, & j = k, \\ 0, & j \neq k, \end{cases} \quad j, k = 1, \dots, n. \quad (5.15)$$

Подставив (5.15) в соотношение (5.14) получаем, учитывая (5.9):

$$\mathbf{Cov}\{y_i, y_l\} = \Sigma \sum_{j=1}^n c_{ij}c_{lj} = \delta_{il}\Sigma = \begin{cases} \Sigma, & i = l, \\ 0, & i \neq l, \end{cases} \quad i, l = 1, \dots, n. \quad (5.16)$$

Таким образом, случайные векторы y_1, \dots, y_n имеют совместное нормальное распределение, некоррелированы, имеют нормальные распределения (5.11) и по теореме 4.4 независимы в совокупности.

Докажем равенство (5.13).

$$\sum_{i=1}^n y_i y_i^{\mathbf{T}} = \sum_{i=1}^n \sum_{j=1}^n c_{ij}x_j \left(\sum_{k=1}^n c_{ik}x_k\right)^{\mathbf{T}} =$$

$$= \sum_{j,k=1}^n \sum_{i=1}^n c_{ij} c_{ik} x_j x_k^T = \sum_{j,k=1}^n \delta_{jk} x_j x_k^T = \sum_{i=1}^n x_i x_i^T.$$

□

Определение 5.1. Говорят, что случайная $(N \times N)$ -матрица A имеет распределение Уишарта $W_N(\Sigma, m)$, если она распределена как матрица

$$A = \sum_{j=1}^m y_j y_j^T,$$

где случайные N -векторы y_1, \dots, y_m независимы в совокупности и одинаково распределены по невырожденному нормальному закону $\mathcal{N}_N(\mathbf{0}_N, \Sigma)$ ($|\Sigma| \neq 0$).

Распределение Уишарта обладает следующими свойствами:

Свойство 1. Если случайные матрицы A_1, \dots, A_K независимы в совокупности и имеют распределения Уишарта:

$$\mathcal{L}\{A_k\} = W_N(\Sigma, m_k), \quad k = 1, \dots, K,$$

то их сумма также имеет распределение Уишарта:

$$\mathcal{L}\left\{\sum_{k=1}^K A_k\right\} = W_N(\Sigma, \sum_{k=1}^K m_k);$$

Свойство 2. Если случайная матрица A имеет распределение Уишарта:

$$\mathcal{L}\{A\} = W_N(\Sigma, m),$$

то $\mathbf{E}\{A\} = m\Sigma$, и

$$\mathcal{L}\{cA\} = W_N(c\Sigma, m),$$

где $c > 0$ – некоторая постоянная величина.

Теорема 5.2 (О распределении оценок максимального правдоподобия параметров многомерного нормального распределения). В условиях теоремы 5.1 МП-оценки

$$\hat{\mu} = \bar{x} \quad \text{и} \quad \hat{\Sigma} = \frac{1}{n} A,$$

определенные в (5.1) и (5.2), независимы и имеют следующие распределения вероятностей:

$$\mathcal{L}\{\bar{x}\} = \mathcal{N}_N\left(\mu, \frac{1}{n}\Sigma\right), \quad (5.17)$$

$$\mathcal{L}\{n\hat{\Sigma}\} = \mathcal{L}\{A\} = W_N(\Sigma, n-1). \quad (5.18)$$

Доказательство. Введем в рассмотрение ортогональную $(n \times n)$ -матрицу C :

$$C = \begin{pmatrix} c_{1,1} & c_{1,2} & \dots & c_{1,n} \\ c_{2,1} & c_{2,2} & \dots & c_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n-1,1} & c_{n-1,2} & \dots & c_{n-1,n} \\ \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \dots & \frac{1}{\sqrt{n}} \end{pmatrix} \quad (5.19)$$

и рассмотрим ортогональные преобразования y_1, \dots, y_n элементов случайной выборки $X = \{x_1, \dots, x_n\}$:

$$y_i = \sum_{j=1}^n c_{ij} \cdot x_j \quad i = 1, \dots, n.$$

Выразим оценки \bar{x} и $\hat{\Sigma}$ через полученные случайные векторы y_1, \dots, y_n .

Согласно лемме 5.1 случайные векторы y_1, \dots, y_n имеют совместное нормальное распределение, некоррелированы, имеют следующие нормальные распределения:

$$\mathcal{L}\{y_i\} = \mathcal{N}_N(\nu_i, \Sigma), \quad i = 1, \dots, n, \quad (5.20)$$

где, с учетом (5.19), ν_i определяются следующим образом:

$$\nu_i = \sum_{j=1}^n c_{ij} \mu = \mu \sum_{j=1}^n c_{ij} \frac{1}{\sqrt{n}} \sqrt{n} = \sqrt{n} \mu \sum_{j=1}^n c_{ij} c_{jn} = \begin{cases} \sqrt{n} \mu, & i = n, \\ 0, & i \neq n, \end{cases}, \quad i = 1, \dots, n. \quad (5.21)$$

Таким образом,

$$\begin{aligned} \mathcal{L}\{y_i\} &= \mathcal{N}_N(\mathbf{0}_N, \Sigma), \quad i = 1, \dots, n-1, \\ \mathcal{L}\{y_n\} &= \mathcal{N}_N(\sqrt{n} \mu, \Sigma). \end{aligned} \quad (5.22)$$

Из (5.21) следует:

$$y_n = \sum_{j=1}^n c_{nj} x_j = \frac{1}{\sqrt{n}} \sum_{j=1}^n x_j = \sqrt{n} \bar{x} \quad (5.23)$$

и

$$\bar{x} = \frac{1}{\sqrt{n}} y_n. \quad (5.24)$$

Из последнего равенства следует, что \bar{x} – линейное преобразование случайного вектора y_n , имеющего невырожденное нормальное распределение $\mathcal{N}_N(\sqrt{n} \mu, \Sigma)$. Следовательно,

$$\mathcal{L}\{\bar{x}\} = \mathcal{N}_N(\mu, \frac{1}{\sqrt{n}} \Sigma). \quad (5.25)$$

Определим закон распределения $\hat{\Sigma}$.

Так как

$$\hat{\Sigma} = \frac{1}{n} A = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{n} \left(\sum_{i=1}^n x_i x_i^T - n \bar{x} \bar{x}^T \right),$$

то с учетом равенства (5.13), получим:

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n-1} y_i y_i^T. \quad (5.26)$$

Согласно (5.23) и определению распределения Уишарта:

$$\mathcal{L}\{n \hat{\Sigma}\} = \mathcal{L}\{A\} = W_N(\Sigma, n-1).$$

Независимость статистических оценок \bar{x} и $\hat{\Sigma}$ непосредственно следует из представлений (5.24) и (5.26), согласно которым \bar{x} и $\hat{\Sigma}$ не содержат общих случайных векторов из y_1, \dots, y_n , так как \bar{x} определяется y_n , а $\hat{\Sigma}$ строится на основе y_1, \dots, y_{n-1} . □

5.3 Вероятностные свойства выборочного среднего и выборочной ковариационной матрицы. Несмещенная выборочная ковариационная матрица

Теорема 5.3. В условиях теоремы 5.1 выборочное среднее \bar{x} является несмещенной оценкой для математического ожидания μ :

$$\mathbf{E}\{\bar{x}\} = \mu,$$

а выборочная ковариационная матрица $\hat{\Sigma}$ — асимптотически несмещенной:

$$\mathbf{E}\{\hat{\Sigma}\} \xrightarrow{n \rightarrow \infty} \Sigma.$$

Несмещенной оценкой ковариационной матрицы является статистика

$$S = \frac{1}{n-1}A, \quad (5.27)$$

где

$$A = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T.$$

Все оценки \bar{x} , $\hat{\Sigma}$ и S сильно состоятельны:

$$\bar{x} \xrightarrow[n \rightarrow \infty]{\mathbf{P}=1} \mu, \quad \hat{\Sigma} \xrightarrow[n \rightarrow \infty]{\mathbf{P}=1} \Sigma, \quad S \xrightarrow[n \rightarrow \infty]{\mathbf{P}=1} \Sigma, \quad n \rightarrow +\infty,$$

Доказательство. Несмещенность выборочного среднего \bar{x} следует из рассуждений:

$$\mathbf{E}\{\bar{x}\} = \mathbf{E}\left\{\frac{1}{n} \sum_{i=1}^n x_i\right\} = \frac{1}{n} \sum_{i=1}^n \mathbf{E}\{x_i\} = \mu.$$

Вычислим $\mathbf{E}\{\hat{\Sigma}\}$ используя равенство (5.26):

$$\mathbf{E}\{\hat{\Sigma}\} = \mathbf{E}\left\{\frac{1}{n} \sum_{i=1}^{n-1} y_i y_i^T\right\} = \frac{1}{n} \sum_{i=1}^{n-1} \mathbf{E}\{y_i y_i^T\} = \frac{n-1}{n} \Sigma. \quad (5.28)$$

Так как

$$\mathbf{E}\{\hat{\Sigma}\} \neq \Sigma,$$

то $\hat{\Sigma}$ — смещенная оценка для Σ ; смещение этой оценки $-b = \frac{1}{n}\Sigma$.

Сильная состоятельность оценок \bar{x} , $\hat{\Sigma}$ и S вытекает из усиленного закона больших чисел для независимых в совокупности одинаково распределенных случайных векторов [11].

Действительно, так как x_1, \dots, x_n — независимы в совокупности и одинаково распределены с математическим ожиданием:

$$\mathbf{E}\{x_k\} = \mu, \quad k = 1, \dots, n,$$

то

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \xrightarrow[n \rightarrow \infty]{\mathbf{P}=1} \mu,$$

а $y_1 y_1^T, \dots, y_{n-1} y_{n-1}^T$ – независимые в совокупности, одинаково распределенные случайные матрицы с математическим ожиданием

$$\mathbf{E}\{y_i y_i^T\} = \Sigma,$$

то

$$\hat{\Sigma} = \frac{n-1}{n} \frac{1}{n-1} \sum_{i=1}^{n-1} y_i y_i^T \xrightarrow[n \rightarrow \infty]{\mathbf{P}=1} \Sigma,$$

и

$$S = \frac{1}{n-1} \sum_{i=1}^{n-1} y_i y_i^T \xrightarrow[n \rightarrow \infty]{\mathbf{P}=1} \Sigma.$$

□

Теорема 5.4. *Элементы выборочной ковариационной матрицы $\hat{\Sigma}$ и несмещенной выборочной ковариационной матрицы S*

$$\hat{\Sigma} = \begin{pmatrix} \hat{\sigma}_{1,1} & \hat{\sigma}_{1,2} & \dots & \hat{\sigma}_{1,N} \\ \hat{\sigma}_{2,1} & \hat{\sigma}_{2,2} & \dots & \hat{\sigma}_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\sigma}_{N,1} & \hat{\sigma}_{N,2} & \dots & \hat{\sigma}_{N,N} \end{pmatrix} \quad u \quad S = \begin{pmatrix} s_{1,1} & s_{1,2} & \dots & s_{1,N} \\ s_{2,1} & s_{2,2} & \dots & s_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ s_{N,1} & s_{N,2} & \dots & s_{N,N} \end{pmatrix} \quad (5.29)$$

имеют совместные асимптотически нормальные распределения с ковариациями, определяемыми элементами истинной ковариационной матрицы Σ :

$$\lim_{n \rightarrow +\infty} n \mathbf{Cov}\{\hat{\sigma}_{kl}, \hat{\sigma}_{rq}\} = \lim_{n \rightarrow +\infty} n \mathbf{Cov}\{s_{kl}, s_{rq}\} = \sigma_{kr} \sigma_{lq} + \sigma_{kq} \sigma_{lr}, \quad k, l, r, q = 1, \dots, N. \quad (5.30)$$

Доказательство. Доказательство теоремы основано на использовании центральной предельной теоремы для одинаково распределенных случайных векторов [1], согласно которой, если $v_1, \dots, v_m, v_{m+1}, \dots \in \mathbb{R}^p$ – последовательность независимых в совокупности, одинаково распределенных случайных p -мерных векторов с математическим ожиданием μ_v и ковариационной матрицей Σ_v :

$$\mu_v = \mathbf{E}\{v_i\}, \quad \Sigma_v = \mathbf{Cov}\{v_i, v_i\}, \quad i = 1, \dots, m, m+1, \dots,$$

то

$$\mathcal{L} \left\{ \frac{1}{\sqrt{m}} \sum_{i=1}^m (v_i - \mu_v) \right\} \xrightarrow{m \rightarrow \infty} \mathcal{N}_p(\mathbf{0}_p, \Sigma_v).$$

Согласно (5.26):

$$n \hat{\Sigma} = (n-1) S = A = \sum_{i=1}^{n-1} y_i y_i^T, \quad (5.31)$$

где случайные N -векторы $y_i = (y_{i1}, \dots, y_{iN})$, независимы в совокупности и, согласно (5.22), распределены по закону $\mathcal{N}_N(\mathbf{0}_N, \Sigma)$.

Из (5.31) следует

$$\mathbf{E} \left\{ \frac{1}{n-1} a_{kl} \right\} = \sigma_{kl}, \quad k, l = 1, \dots, N. \quad (5.32)$$

Введем в рассмотрение нормированные элементы a_{kl}° матрицы A :

$$a_{kl}^\circ = \frac{1}{\sqrt{n-1}} (a_{kl} - (n-1) \sigma_{kl}), \quad k, l = 1, \dots, N, \quad (5.33)$$

из которых, учитывая, что по построению $a_{kl}^\circ = a_{lk}^\circ$, $k, l = 1, \dots, N$, составим $\frac{N(N+1)}{2}$ -вектор

$$a^\circ = (a_{11}^\circ, a_{21}^\circ, a_{22}^\circ, a_{31}^\circ, a_{32}^\circ, a_{33}^\circ, a_{41}^\circ, \dots, a_{NN}^\circ)^\mathbf{T}$$

который имеет асимптотически нормальное $(n \rightarrow +\infty)$ $\frac{N(N+1)}{2}$ -мерное распределение с нулевыми математическими ожиданиями и ковариациями:

$$\begin{aligned} \mathbf{Cov}\{a_{kl}^\circ, a_{rq}^\circ\} &= \mathbf{Cov}\left\{\frac{1}{\sqrt{n-1}}(a_{kl} - (n-1)\sigma_{kl}), \frac{1}{\sqrt{n-1}}(a_{rq} - (n-1)\sigma_{rq})\right\} = \\ &= \mathbf{Cov}\{y_{ik}y_{il}, y_{ir}y_{iq}\} = \mathbf{E}\{y_{ik}y_{il}y_{ir}y_{iq}\} - \mathbf{E}\{y_{ik}y_{il}\}\mathbf{E}\{y_{ir}y_{iq}\} = \\ &= \sigma_{kl}\sigma_{rq} + \sigma_{kr}\sigma_{lq} + \sigma_{kq}\sigma_{lr} - \sigma_{kl}\sigma_{rq} = \sigma_{kr}\sigma_{lq} + \sigma_{kq}\sigma_{lr}, \quad k, l, r, q = 1, \dots, N, \end{aligned}$$

где использована известная формула (4.4) для моментов четвертого порядка многомерного нормального распределения.

Очевидное соотношение:

$$\begin{aligned} \lim_{n \rightarrow +\infty} n \mathbf{Cov}\{\hat{\sigma}_{kl}, \hat{\sigma}_{rq}\} &= \lim_{n \rightarrow +\infty} n \mathbf{Cov}\{s_{kl}, s_{rq}\} = \\ &= \lim_{n \rightarrow +\infty} \mathbf{Cov}\left\{\frac{1}{\sqrt{n-1}}(a_{kl} - (n-1)\sigma_{kl}), \frac{1}{\sqrt{n-1}}(a_{rq} - (n-1)\sigma_{rq})\right\}, \quad k, l, r, q = 1, \dots, N, \end{aligned}$$

завершает доказательство. □

6 Корреляционный анализ (статистическое исследование зависимостей)

6.1 Выборочный коэффициент корреляции и его свойства

6.1.1 Определение и основные свойства выборочного коэффициента корреляции

Пусть в пространстве из N , ($N \geq 1$) признаков наблюдается выборка $X = \{x_1, \dots, x_n\}$ объема n , образованная наблюдениями с вектором математического ожидания μ и ковариационной $(N \times N)$ -матрицей Σ :

$$\mu = \mathbf{E}\{x_i\} \in \mathbb{R}^N \text{ и } \Sigma = \mathbf{Cov}\{x_i, x_i\}, i = 1, \dots, n.$$

Элементы σ_{kl} , $k, l = 1, \dots, N$, ковариационной матрицы Σ , описывают *парные зависимости* признаков между собой. Однако на практике для исследования парной зависимости признаков чаще используется коэффициент корреляции, называемый еще *парной корреляцией*:

$$\rho_{kl} = \frac{\sigma_{kl}}{\sqrt{\sigma_{kk}\sigma_{ll}}}, k, l = 1, \dots, N, \quad (6.1)$$

В тех случаях, когда вектор математического ожидания μ и ковариационная матрица Σ неизвестны, для исследования зависимостей используют статистическую оценку парного коэффициента корреляции, вычисленную по выборке X объема n – *выборочный коэффициент корреляции*.

Определение 6.1. Выборочным коэффициентом корреляции $\hat{\rho}_{kl}$ между компонентами \tilde{x}_k и \tilde{x}_l случайного N -мерного вектора $x = (\tilde{x}_1, \dots, \tilde{x}_N)^T \in \mathbb{R}^N$, вычисленным по случайной выборке $X = \{x_1, \dots, x_n\}$ объема n , называется статистика

$$\hat{\rho}_{kl} = r_{kl} = \frac{\hat{\sigma}_{kl}}{\sqrt{\hat{\sigma}_{kk}\hat{\sigma}_{ll}}}, k, l = 1, \dots, N, \quad (6.2)$$

где $\hat{\sigma}_{kl}$, $\hat{\sigma}_{kk}$, $\hat{\sigma}_{ll}$ – элементы выборочной ковариационной матрицы $\hat{\Sigma}$, $k, l = 1, \dots, N$.

Так как

$$\hat{\sigma}_{kl} = \frac{1}{n}a_{kl}, k, l = 1, \dots, N,$$

где статистики a_{kl} , определены равенствами

$$a_{kl} = \sum_{i=1}^n (x_{ik} - \bar{x}_k)(x_{il} - \bar{x}_l), k, l = 1, \dots, N, \quad (6.3)$$

и являются элементами матрицы A :

$$A = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T,$$

то эквивалентным соотношению (6.2) является равенство

$$\hat{\rho}_{kl} = r_{kl} = \frac{a_{kl}}{\sqrt{a_{kk}a_{ll}}}, k, l = 1, \dots, N. \quad (6.4)$$

Свойства выборочного коэффициента корреляции $\hat{\rho}_{kl}$

Свойство 1. $-1 \leq \hat{\rho}_{kl} \leq 1$.

Свойство 2. $\hat{\rho}_{kl}$ инвариантен относительно масштабного преобразования признаков в выборке, то есть не изменяют своих значений при умножении признаков в наблюдениях $x_i = (x_{i1}, \dots, x_{iN})^T$ на соответствующие константы: $x_i := (b_1 x_{i1}, \dots, b_N x_{iN})^T$, $i = 1, \dots, n$.

Свойство 3. Вычисленный по случайной выборке X объема n выборочный коэффициент корреляции $\hat{\rho}_{kl}$ – строго состоятельная оценка парного коэффициента корреляции ρ_{kl} :

$$\hat{\rho}_{kl} \xrightarrow[n \rightarrow \infty]{P=1} \rho_{kl}. \quad (6.5)$$

Предположим теперь, что наблюдения из выборки $X = \{x_1, \dots, x_n\}$ независимы в совокупности и имеют невырожденное многомерное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$). Исследуем дополнительные свойства выборочного коэффициента корреляции (6.2), который в этом случае является оценкой максимального правдоподобия для коэффициента корреляции (6.1) (поскольку (6.2) – функциональное преобразование выборочной ковариационной матрицы, являющейся оценкой максимального правдоподобия).

Теорема 6.1 (О распределении выборочного коэффициента корреляции). [1] Пусть \tilde{x}_k и \tilde{x}_l – компоненты случайного вектора x , распределенного по закону $\mathcal{N}_N(\mu, \Sigma)$, причем $\rho_{ij} = \text{Corr}\{x_i, x_j\} = 0$. Тогда выборочный коэффициент корреляции, вычисленный по выборке $X = \{x_1, \dots, x_n\}$, имеет плотность распределения вероятностей $p_0(r)$ вида

$$p_0(r) = \frac{\Gamma\left(\frac{n-1}{2}\right)}{\sqrt{\pi}\Gamma\left(\frac{n-2}{2}\right)} (1-r^2)^{\frac{n-4}{2}}, \quad (6.6)$$

где

$$\Gamma(a) = \int_0^\infty e^{-u} u^{a-1} du, \quad \Gamma(0, 5) = \sqrt{\pi}, \quad |r| < 1.$$

Очевидно, что функция $p_0(r)$ четная, то есть

$$p_0(r) = p_0(-r). \quad (6.7)$$

Из (6.7) следует, что для функции распределения вероятностей $F_0(r)$:

$$F_0(r) = \int_{-\infty}^r p_0(u) du, \quad (6.8)$$

будет выполняться равенство

$$F_0(r) = 1 - F_0(-r). \quad (6.9)$$

6.1.2 Проверка гипотезы о значимости коэффициента корреляции

Пусть имеется выборка $X = \{x_1, \dots, x_n\}$ из распределения вероятностей $\mathcal{N}_N(\mu, \Sigma)$. Необходимо при заданном уровне значимости α проверить гипотезу о равенстве нулю коэффициента корреляции ρ_{ij} :

$$\begin{aligned} H_0 & : \rho_{ij} = 0; \\ H_1 = \overline{H_0} & : \rho_{ij} \neq 0. \end{aligned} \quad (6.10)$$

Решение задачи (6.10) формулируется при помощи статистики r_{ij} следующим образом:

$$\begin{aligned} H_0 & : |r_{ij}| < \Delta; \\ H_1 = \overline{H_0} & : |r_{ij}| \geq \Delta, \end{aligned} \quad (6.11)$$

где Δ – порог критерия, который определяется следующим образом:

$$\begin{aligned}\alpha &= \mathbf{P}\{H_1 | H_0\} = \mathbf{P}\{|r_{ij}| \geq \Delta | H_0\} = \mathbf{P}\{r_{ij} \leq -\Delta | H_0\} + \mathbf{P}\{r_{ij} \geq \Delta | H_0\} = \\ &= \mathbf{P}\{r_{ij} \leq -\Delta | H_0\} + 1 - \mathbf{P}\{r_{ij} < \Delta | H_0\} = F_0(-\Delta) + 1 - F_0(\Delta),\end{aligned}\quad (6.12)$$

где $F_0(\Delta)$ определена в (6.8). Учитывая равенство (6.9) можно записать

$$\alpha = 2 \cdot (1 - F_0(\Delta)),$$

из которого следует, что порог критерия (6.12) определяется равенством

$$\Delta = F_0^{-1}\left(1 - \frac{\alpha}{2}\right).$$

6.1.3 Z-преобразование и Z-статистика Фишера

Теорема 6.2 (Об асимптотической нормальности выборочного коэффициента корреляции). [12] В условиях теоремы 5.1 и при условии, что соответствующее истинное абсолютное значение коэффициента корреляции ρ_{ij} не равно единице:

$$|\rho_{ij}| \neq 1,$$

выборочный коэффициент корреляции r_{ij} из (6.2) имеет асимптотически нормальное распределение:

$$\mathcal{L}\{\sqrt{n-1}(r_{ij} - \rho_{ij})\} \xrightarrow{n \rightarrow \infty} \mathcal{N}_1(0, (1 - \rho_{ij}^2)^2), \quad ; i, j = 1, \dots, N. \quad (6.13)$$

Теорема 6.3 (Об асимптотической нормальности функционально преобразованной последовательности). [1] Пусть u_1, \dots, u_m – асимптотически нормальная сходящаяся по вероятности случайная последовательность:

$$u_m \xrightarrow[m \rightarrow \infty]{\mathbf{P}} b,$$

$$\mathcal{L}\{\sqrt{m}(u_m - b)\} \xrightarrow{m \rightarrow \infty} \mathcal{N}_1(0, \sigma^2),$$

где $b \in \mathbb{R}$, а $\sigma^2 > 0$.

Если функция

$$f = f(z) : \mathbb{R}^1 \rightarrow \mathbb{R}^1$$

имеет в окрестности точки $z = b$ первую и вторую производные и

$$\phi_b = f'(z)|_{z=b} \neq 0, \quad (6.14)$$

то

$$\mathcal{L}\{\sqrt{m}(f(u_m) - f(b))\} \xrightarrow{m \rightarrow \infty} \mathcal{N}_1(0, \beta), \quad (6.15)$$

где

$$\beta = \sigma^2 \cdot \phi_b^2. \quad (6.16)$$

Доказательство. Доказательство основано на разложении функции $f(z)$ в точке $z = b$ в ряд Тейлора с остаточным членом в форме Лагранжа и на вычислении соответствующих моментов. \square

Следствие 6.1. Пусть $\Psi = \Psi(v) : \mathbb{R}^1 \rightarrow \mathbb{R}^1$ – дважды дифференцируема в точке $v = \rho_{ij}$, и такая, что

$$\frac{d}{dv}\Psi(v)|_{v=\rho_{ij}} \neq 0,$$

тогда в условиях теоремы 6.2 для выборочного коэффициента корреляции r_{ij} выполняется ($n \rightarrow +\infty$):

$$\mathcal{L}\{\sqrt{n-1}(\Psi(r_{ij}) - \Psi(\rho_{ij}))\} \rightarrow \mathcal{N}_1\left(0, \left(\frac{d}{dv}\Psi(v)\Big|_{v=\rho_{ij}}\right)^2 (1 - \rho_{ij}^2)^2\right). \quad (6.17)$$

Доказательство. Доказательство непосредственно следует из результата теоремы 6.2 и теоремы 6.3. \square

Замечание 6.1.1. В соотношениях (6.13) и (6.17) вместо $n-1$ можно использовать n (поскольку $n \rightarrow +\infty$), однако считается, что $n-1$ «точнее».

Выберем преобразование $\Psi = \Psi(v) : \mathbb{R}^1 \rightarrow \mathbb{R}^1$ так, чтобы дисперсия предельного нормального распределения в (6.17) была единичной:

$$\left(\frac{d}{dv}\Psi(v)\Big|_{v=\rho_{ij}}\right)^2 (1 - \rho_{ij}^2)^2 = 1.$$

Получим дифференциальное уравнение:

$$\frac{d}{dv}\Psi(v)\Big|_{v=\rho_{ij}} = \frac{1}{1 - \rho_{ij}^2} = \frac{1}{(1 - \rho_{ij})(1 + \rho_{ij})} = \frac{1}{2} \left(\frac{1}{1 + \rho_{ij}} + \frac{1}{1 - \rho_{ij}} \right),$$

решая которое, найдем

$$\Psi(v) = \frac{1}{2} \ln \frac{1+v}{1-v}.$$

Определение 6.2. Пусть r_{ij} – выборочный коэффициент корреляции (6.2), тогда статистика

$$Z = Z(r_{ij}) = \Psi(r_{ij}) = \frac{1}{2} \ln \frac{1+r_{ij}}{1-r_{ij}} \quad (6.18)$$

называется Z -статистикой Фишера, а преобразование $\Psi(\cdot)$ – Z -преобразованием Фишера.

Следствие 6.2. В условиях теоремы 6.2 случайная величина $\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}))$ имеет при $n \rightarrow +\infty$ стандартное нормальное распределение $\mathcal{N}_1(0, 1)$.

6.1.4 Проверка гипотезы о значении коэффициента корреляции

Z -статистика Фишера (6.18) имеет большое прикладное значение – с ее помощью проверяются гипотезы о значении коэффициента корреляции:

$$\begin{aligned} H_0 & : \rho_{ij} = \rho_{ij}^o; \\ H_1 = \overline{H_0} & : \rho_{ij} \neq \rho_{ij}^o, \end{aligned} \quad (6.19)$$

где ρ_{ij}^o – предполагаемое значение коэффициента корреляции. Очевидно, что задача (6.19) эквивалентна следующей задаче

$$\begin{aligned} H_0 & : \sqrt{n-1}(\rho_{ij} - \rho_{ij}^o) = 0; \\ H_1 = \overline{H_0} & : \sqrt{n-1}(\rho_{ij} - \rho_{ij}^o) \neq 0. \end{aligned} \quad (6.20)$$

Учитывая следствие 6.2, а именно тот факт, что при $n \rightarrow +\infty$ и верной гипотезе H_0

$$\mathcal{L}\{\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o))\} = \mathcal{N}_1(0, 1),$$

получим при следующий критерий для проверки гипотез H_0, H_1 :

$$\begin{cases} H_0 & : |\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o))| \leq \Delta; \\ H_1 = \overline{H_0} & : |\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o))| > \Delta, \end{cases} \quad (6.21)$$

где порог критерия Δ определяется по наперед заданному малому значению уровня значимости α (вероятность принять гипотезу H_1 при условии, что верна гипотеза H_0):

$$\alpha = \mathbf{P}\{H_1|H_0\} = 1 - \mathbf{P}\{H_0|H_0\} \in (0, 1),$$

$$\begin{aligned} \alpha &= \mathbf{P}\{H_1 | H_0\} = \mathbf{P}\{|\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o))| > \Delta | H_0\} = \\ &= \mathbf{P}\{\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o)) \leq -\Delta | H_0\} + \mathbf{P}\{\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o)) \geq \Delta | H_0\} = \\ &= \mathbf{P}\{\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o)) \leq -\Delta | H_0\} + 1 - \mathbf{P}\{\sqrt{n-1}(Z(r_{ij}) - Z(\rho_{ij}^o)) < \Delta | H_0\} = \\ &= \Phi(-\Delta) + 1 - \Phi(\Delta) = 2(1 - \Phi(\Delta)), \end{aligned} \quad (6.22)$$

где $\Phi(\cdot)$ – функция распределения стандартного нормального закона распределения. Таким образом:

$$\Delta = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right), \quad (6.23)$$

то есть Δ – квантиль уровня $1 - \frac{\alpha}{2}$ стандартного нормального закона $\mathcal{N}_1(0, 1)$.

Замечание 6.2.1. Построенный критерий (6.21), (6.23) позволяет также определить доверительный интервал [12] уровня $1 - \alpha$ для ρ_{ij} :

$$\mathbf{P}\left\{\text{th}\left(Z(r_{ij}) - \frac{\Delta}{\sqrt{n-1}}\right) < \rho_{ij} < \text{th}\left(Z(r_{ij}) + \frac{\Delta}{\sqrt{n-1}}\right)\right\} = 1 - \alpha,$$

где $\Psi^{-1}(v)$ – преобразование, обратное к преобразованию Фишера $\Psi(v)$:

$$\Psi^{-1}(v) = \frac{e^v - e^{-v}}{e^v + e^{-v}} = \text{th}(v),$$

$$\Psi(v) = \frac{1}{2} \ln \frac{1+v}{1-v}.$$

Замечание 6.2.2. На практике чаще всего проверяются гипотезы:

$$\begin{aligned} H_0 & : \rho_{ij} = 0; \\ H_1 = \overline{H_0} & : \rho_{ij} \neq 0. \end{aligned}$$

Если принимается H_0 , то i -й и j -й признаки ($k \neq l = \overline{1, N}$) считаются некоррелированными, а в силу предположения их совместной нормальности – и независимыми. Поэтому гипотезу H_0 иногда называют гипотезой независимости. Проверяются гипотезы H_0, H_1 при помощи критерия (6.21), (6.23) ($\rho_{ij}^o = 0$, и $Z(\rho_{ij}^o) = Z(0) = 0$).

В случае выборки «малого объема» (асимптотика $n \rightarrow +\infty$ не предполагается) для проверки гипотез H_0, H_1 считается целесообразным использование другого критерия:

$$\begin{cases} H_0 & : \sqrt{n-2} \frac{|r_{ij}|}{\sqrt{1-r_{ij}}} \leq \Delta; \\ H_1 = \overline{H_0} & : \sqrt{n-2} \frac{|r_{ij}|}{\sqrt{1-r_{ij}}} > \Delta, \end{cases}$$

где относительно статистики критерия в пункте 6.2 (теорема 6.5, соотношение (6.29)) будет установлено, что при истинной гипотезе H_0 она является модулем случайной величины, имеющей t -распределение Стьюдента с $n-2$ степенями свободы, и порог критерия $\Delta = F_{t_{n-2}}^{-1}(1 - \frac{\alpha}{2})$ – соответствующая квантиль уровня $1 - \frac{\alpha}{2}$ ($\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$). Однако данный критерий «более чувствителен» к отклонениям распределения вероятностей выборочных значений компонент от нормального закона, и при «достаточно большом» объеме выборки n лучше использовать критерий (6.21), (6.23), основанный на Z -статистике Фишера (6.18).

6.2 Выборочный частный коэффициент корреляции и его свойства

6.2.1 Определение и свойства выборочного частного коэффициента корреляции

Пусть случайный N -мерный вектор признаков $x = (\tilde{x}_1, \dots, \tilde{x}_N) \in \mathbb{R}^N$ имеет невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$ с неизвестным математическим ожиданием μ и неизвестной ковариационной матрицей Σ , а $X = \{x_1, \dots, x_n\}$ – случайная выборка объема n , $n > N$ наблюдений за x .

Оценим по выборке X частный коэффициент корреляции (4.29). Согласно разбиению случайного вектора

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix},$$

на подвекторы $x^{(1)} \in \mathbb{R}^m$ и $x^{(2)} \in \mathbb{R}^{N-m}$ ($m < N$), где

$$x^{(1)} = (\tilde{x}_1, \dots, \tilde{x}_m)^{\mathbf{T}} \in \mathbb{R}^m, \quad x^{(2)} = (\tilde{x}_{m+1}, \dots, \tilde{x}_N)^{\mathbf{T}} \in \mathbb{R}^{N-m},$$

порождающего разбиения N -вектора математического ожидания $\mu \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицы Σ – на соответствующие блоки:

$$\mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}, \quad \Sigma_{21} = \Sigma_{12}^{\mathbf{T}}. \quad (6.24)$$

построим соответствующие разбиения их статистических оценок – арифметического среднего \bar{x} и выборочной ковариационной матрицы $\hat{\Sigma}$:

$$\bar{x} = \begin{pmatrix} \bar{x}^{(1)} \\ \bar{x}^{(2)} \end{pmatrix}, \quad \hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{11} & \hat{\Sigma}_{12} \\ \hat{\Sigma}_{21} & \hat{\Sigma}_{22} \end{pmatrix} = \frac{1}{n} \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}, \quad A_{21} = A_{12}^{\mathbf{T}}, \quad (6.25)$$

где $(N \times N)$ -матрица A определена в (5.3).

Введем в рассмотрение $(m \times m)$ -матрицу $A_{11|2}$:

$$A_{11|2} = A_{11} - A_{12}A_{22}^{-1}A_{12}^{\mathbf{T}},$$

и обозначим (k, l) -й элемент $A_{11|2}$ через $a_{kl|m+1, \dots, N}$:

$$a_{kl|m+1, \dots, N} = a_{kl} - a_{(k)}^{\mathbf{T}} A_{22}^{-1} a_{(l)}, \quad k, l = 1, \dots, m,$$

где $a_{(k)}^{\mathbf{T}}$ – k -я строка матрицы A_{12} , $k = 1, \dots, m$.

Теорема 6.4. Пусть элементы случайной выборки $X = \{x_1, \dots, x_n\}$ объема $n > N$ независимы в совокупности и имеют невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$, тогда оценкой максимального правдоподобия для частного коэффициента корреляции (4.29) является статистика:

$$r_{kl|m+1, \dots, N} = \frac{a_{kl|m+1, \dots, N}}{\sqrt{a_{kk|m+1, \dots, N} a_{ll|m+1, \dots, N}}}, \quad k, l = 1, \dots, m. \quad (6.26)$$

Доказательство. Доказательство очевидно и основано на подстановке в (4.29) с учетом обозначений (6.25) вместо неизвестной ковариационной матрицы Σ ее МП-оценки

$$\hat{\Sigma} = \frac{1}{n} A$$

по выборке X , полученной в теореме 5.1, условия которой здесь выполняются. \square

Определение 6.3. Выборочным частным коэффициентом корреляции между k -й и l -й компонентами N -мерного случайного вектора $x \in \mathbb{R}^N$ при условии, что компоненты $\tilde{x}_{m+1}, \dots, \tilde{x}_N$ фиксированы, вычисленным по случайной выборке $X = \{x_1, \dots, x_n\}$ объема $n > N$ называется статистика

$$\hat{\rho}_{kl|m+1, \dots, N} = r_{kl|m+1, \dots, N} = \frac{a_{kl|m+1, \dots, N}}{\sqrt{a_{kk|m+1, \dots, N} a_{ll|m+1, \dots, N}}}, \quad k, l = 1, \dots, m. \quad (6.27)$$

Свойства выборочного частного коэффициента корреляции

Свойство 1. $-1 \leq \hat{\rho}_{kl|m+1, \dots, N} \leq 1$.

Свойство 2. Вычисленный по случайной выборке X объема n выборочный частный коэффициент корреляции $\hat{\rho}_{kl|m+1, \dots, N}$ – состоятельная оценка парного коэффициента корреляции $\rho_{kl|m+1, \dots, N}$:

$$\hat{\rho}_{kl|m+1, \dots, N} \xrightarrow[n \rightarrow \infty]{\mathbf{P}} \rho_{kl|m+1, \dots, N}. \quad (6.28)$$

Найдем распределение вероятностей случайной матрицы

$$A_{11|2} = A_{11} - A_{12} A_{22}^{-1} A_{12}^T,$$

элементы которой служат для определения выборочного частного (6.27) коэффициента корреляции.

Теорема 6.5. В условиях теоремы 6.4 справедливы следующие утверждения:

1) случайная матрица $A_{11|2}$:

$$A_{11|2} = A_{11} - A_{12} A_{22}^{-1} A_{12}^T$$

имеет распределение Уишарта

$$\mathcal{L}\{A_{11|2}\} = W_m(\Sigma_{11|2}, n - (N - m) - 1)$$

и не зависит от матрицы $A_{12} A_{22}^{-1} A_{12}^T$, которая при $\Sigma_{12} = \mathbf{0}_{m \times (N-m)}$ также имеет распределение Уишарта:

$$\mathcal{L}\{A_{12} A_{22}^{-1} A_{12}^T | \Sigma_{12} = \mathbf{0}_{m \times (N-m)}\} = W_m(\Sigma_{11|2}, N - m).$$

2) выборочный частный коэффициент корреляции $r_{kl|m+1, \dots, N}$ ($k, l = \overline{1, m}$), вычисленный по выборке объема n , распределен так же, как соответствующий ему обычный выборочный коэффициент корреляции r_{kl} , подсчитанный по выборке объема $n - (N - m)$ с истинным значением коэффициента корреляции $\rho_{kl} := \rho_{kl|m+1, \dots, N}$, $k, l = \overline{1, m}$.

3) для выборочного коэффициента корреляции r_{kl} ($k, l = \overline{1, N}$), определенного по выборке объема n , при нулевом истинном значении: $\rho_{kl} = 0$, статистика

$$\sqrt{n-2} \frac{|r_{kl}|}{\sqrt{1-r_{kl}^2}} \quad (6.29)$$

распределена как модуль случайной величины, имеющей t -распределение Стьюдента с $n-2$ степенями свободы.

6.2.2 Проверка гипотезы о значении частного коэффициента корреляции

Второе утверждение теоремы 6.5 позволяет использовать для проверки гипотез о значении частного коэффициента корреляции

$$\begin{aligned} H_0 &: \rho_{kl|m+1, \dots, N} = \rho_{kl|m+1, \dots, N}^o; \\ H_1 = \overline{H_0} &: \rho_{kl|m+1, \dots, N} \neq \rho_{kl|m+1, \dots, N}^o, \end{aligned} \quad (6.30)$$

где $\rho_{kl|m+1, \dots, N}^o$ – предполагаемое значение частного коэффициента корреляции, все критерии из параграфа 6.1, заменив в них n на $n - (N - m)$.

Гипотеза независимости:

$$\rho_{kl|m+1, \dots, N}^o = 0,$$

в данном случае называется *гипотезой условной независимости*, $k, l = 1, \dots, m$.

6.3 Выборочный множественный коэффициент корреляции и его свойства

Определение 6.4. Выборочным множественным коэффициентом корреляции *между* между компонентой \tilde{x}_k и компонентами подвектора $x^{(2)} = (\tilde{x}_{m+1}, \dots, \tilde{x}_N)^T$, вычисленным по случайной выборке $X = \{x_1, \dots, x_n\}$ объема $n > N$ называется статистика

$$R_{k,m+1, \dots, N} = \sqrt{\frac{a_{(k)}^T A_{22}^{-1} a_{(k)}}{a_{kk}}} \quad k = 1, \dots, m, \quad (6.31)$$

где $a_{(k)}^T$ – k -я строка матрицы A_{12} , $k = 1, \dots, m$, являющейся блоком разбиения определенной в (5.3) $(N \times N)$ -матрицы A :

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}, \quad A_{21} = A_{12}^T.$$

Теорема 6.6. В условиях теоремы 6.4 статистика

$$\frac{n - (N - m) - 1}{N - m} \frac{R_{k,m+1, \dots, N}^2}{1 - R_{k,m+1, \dots, N}^2}$$

при нулевом истинном значении соответствующего множественного коэффициента корреляции:

$$\overline{R}_{k,m+1, \dots, N} = 0,$$

имеет F -распределение Фишера с $N - m$ и $n - (N - m) - 1$ степенями свободы:

$$\mathcal{L} \left\{ \frac{n - (N - m) - 1}{N - m} \frac{R_{k,m+1, \dots, N}^2}{1 - R_{k,m+1, \dots, N}^2} \middle| \overline{R}_{k,m+1, \dots, N} = 0 \right\} = F_{N-m, n-(N-m)-1}, \quad k = 1, \dots, m. \quad (6.32)$$

Воспользуемся выборочным множественным коэффициентом корреляции (6.31) для проверки гипотез о некоррелированности k -го признака ($1 \leq k \leq m$) с группой из $N - m$ выбранных признаков:

$$\begin{aligned} H_0 & : \bar{R}_{k,m+1,\dots,N} = 0; \\ H_1 = \bar{H}_0 & : \bar{R}_{k,m+1,\dots,N} \neq 0. \end{aligned}$$

Гипотеза H_0 в предположении нормальности также называется *гипотезой независимости*.

Учтем, что статистика в (6.32) монотонно убывает с уменьшением $R_{k,m+1,\dots,N}^2$ (приближением его значения к нулю), и построим следующий критерий для проверки гипотез H_0, H_1 :

$$\begin{cases} H_0 & : \frac{n-(N-m)-1}{N-m} \frac{R_{k,m+1,\dots,N}^2}{1-R_{k,m+1,\dots,N}^2} \leq \Delta; \\ H_1 = \bar{H}_0 & : \frac{n-(N-m)-1}{N-m} \frac{R_{k,m+1,\dots,N}^2}{1-R_{k,m+1,\dots,N}^2} > \Delta, \end{cases} \quad (6.33)$$

где порог критерия Δ определяется с учетом (6.32) по заданному малому значению уровня значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$:

$$\Delta = F_{N-m, n-(N-m)-1}^{-1}(1 - \alpha) \quad (6.34)$$

квантиль уровня $1 - \alpha$ от F -распределения с $N - m$ и $n - (N - m) - 1$ степенями свободы.

Замечание 6.4.1. Критерий (6.33), (6.34) в множественном регрессионном анализе называется *проверкой регрессии на значимость*. Если гипотеза H_0 отвергается, то регрессия k -го признака ($1 \leq k \leq m$) на выбранные $N - m$ признаков считается значимой, и возможно прогнозирование значения \tilde{x}_k этого признака по $x^{(2)} \in \mathbb{R}^{N-m}$ при помощи подстановочной статистической оценки функции регрессии (4.17), (4.42):

$$\hat{\tilde{x}}_k = \bar{x}_k + a_{(k)}^{\mathbf{T}} A_{22}^{-1} (x^{(2)} - \bar{x}^{(2)}). \quad (6.35)$$

6.4 Проверка общих гипотез о независимости

Пусть случайный N -мерный вектор $x \in \mathbb{R}^N$, имеющий невырожденное нормальное распределение $\mathcal{N}_N(\mu, \Sigma)$, разбит на $2 \leq K \leq N$ подвекторов-блоков:

$$x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(K)} \end{pmatrix}; \quad x^{(k)} \in \mathbb{R}^{N^{(k)}}, \quad k = 1, \dots, K; \quad N^{(1)} + \dots + N^{(K)} = N.$$

Данное разбиение инициирует соответствующие разбиения N -мерного вектора математического ожидания $\mu \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицы Σ :

$$\mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \\ \vdots \\ \mu^{(K)} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \dots & \Sigma_{1K} \\ \Sigma_{21} & \Sigma_{22} & \dots & \Sigma_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{K1} & \Sigma_{K2} & \dots & \Sigma_{KK} \end{pmatrix}.$$

По случайной выборке $X = \{x_1, \dots, x_n\}$ объема $n > N$, образованной независимыми в совокупности наблюдениями над случайным вектором $x \in \mathbb{R}^N$, необходимо проверить гипотезу о том, что подвекторы (группы признаков) $\{x^{(k)} \in \mathbb{R}^{N^{(k)}}, \quad k = 1, \dots, K\}$ независимы в совокупности между собой.

В силу того что подвекторы $\{x^{(k)} \in \mathbb{R}^{N^{(k)}}, k = 1, \dots, K\}$ имеют совместное нормальное распределение, их независимость:

$$n_N(x|\mu, \Sigma) = \prod_{k=1}^K n_{N^{(k)}}(x^{(k)}|\mu^{(k)}, \Sigma_{kk}),$$

согласно теореме 4.4 эквивалентна их некоррелированности:

$$\mathbf{Cov}\{x^{(k)}, x^{(j)}\} = \Sigma_{kj} = \mathbf{0}_{N^{(k)} \times N^{(j)}}, \quad k \neq j, \quad k, j = 1, \dots, K,$$

и гипотеза независимости сводится к гипотезе

$$H_0 : \Sigma_{kj} = \delta_{kj} \Sigma_{kk}, \quad k, j = 1, \dots, K, \quad (6.36)$$

то есть к гипотезе о том, что Σ имеет вид:

$$\Sigma_0 = \begin{pmatrix} \Sigma_{11} & 0 & \dots & 0 \\ 0 & \Sigma_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Sigma_{KK} \end{pmatrix},$$

при альтернативе общего вида: $H_1 = \overline{H_0}$.

Воспользуемся для проверки сложных гипотез H_0, H_1 критерием отношения правдоподобия [12]. Определим обобщенную статистику отношения правдоподобия:

$$\lambda = \frac{\max_{(\mu, \Sigma) \in \Theta_0} L_N(\mu, \Sigma)}{\max_{(\mu, \Sigma) \in \Theta} L_N(\mu, \Sigma)} \in [0, 1], \quad (6.37)$$

где $L_N(\mu, \Sigma)$ – функция правдоподобия, вычисленная по выборке X объема n :

$$L_N(\mu, \Sigma) = \prod_{i=1}^n n_N(x_i|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{1}{2}nN} |\Sigma|^{\frac{1}{2}n}} \exp \left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right); \quad (6.38)$$

где Θ – множество всех допустимых значений параметров:

$$\Theta = \{(\mu, \Sigma) : \mu \in \mathbb{R}^N, \Sigma = \Sigma^T \succ 0\},$$

а Θ_0 – множество значений параметров, соответствующее гипотезе H_0 :

$$\Theta_0 = \{(\mu, \Sigma) : \mu \in \mathbb{R}^N, \Sigma = \text{diag}\{\Sigma_{11}, \dots, \Sigma_{KK}\}, \Sigma_{kk} = \Sigma_{kk}^T \succ 0, k = 1, \dots, K\}.$$

Критерий отношения правдоподобия, основанный на статистике (6.37), имеет вид

$$\begin{cases} H_0 & : \lambda \geq \lambda_o(\alpha); \\ H_1 = \overline{H_0} & : \lambda < \lambda_o(\alpha), \end{cases} \quad (6.39)$$

где порог критерия $\lambda_o = \lambda_o(\alpha)$ определяется по наперед заданному малому значению уровня значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$.

Разобьем вычисленные по выборке X объема $n > N$ выборочное среднее (5.1):

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

выборочную ковариационную матрицу (5.2):

$$\hat{\Sigma} = \frac{1}{n}A, \quad A = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T = (a_{kl})_{k,l=1}^N,$$

а также выборочную корреляционную матрицу:

$$\mathcal{R} = (r_{kl})_{k,l=1}^N, \quad r_{kl} = \frac{a_{kl}}{\sqrt{a_{kk}a_{ll}}}, \quad k, l = \overline{1, N},$$

на соответствующие блоки:

$$\bar{x} = \begin{pmatrix} \bar{x}^{(1)} \\ \bar{x}^{(2)} \\ \vdots \\ \bar{x}^{(K)} \end{pmatrix}, \quad \hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{11} & \hat{\Sigma}_{12} & \dots & \hat{\Sigma}_{1K} \\ \hat{\Sigma}_{21} & \hat{\Sigma}_{22} & \dots & \hat{\Sigma}_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\Sigma}_{K1} & \hat{\Sigma}_{K2} & \dots & \hat{\Sigma}_{KK} \end{pmatrix}; \quad \Sigma_{kj} = \frac{1}{n}A_{kj}, \quad k, j = 1, \dots, K,$$

$$A = \begin{pmatrix} A_{11} & A_{12} & \dots & A_{1K} \\ A_{21} & A_{22} & \dots & A_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ A_{K1} & A_{K2} & \dots & A_{KK} \end{pmatrix}; \quad \mathcal{R} = \begin{pmatrix} \mathcal{R}_{11} & \mathcal{R}_{12} & \dots & \mathcal{R}_{1K} \\ \mathcal{R}_{21} & \mathcal{R}_{22} & \dots & \mathcal{R}_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{R}_{K1} & \mathcal{R}_{K2} & \dots & \mathcal{R}_{KK} \end{pmatrix}.$$

Теорема 6.7. Пусть выборка $X = \{x_1, \dots, x_i\}$ объема $n > N$ образована независимыми в совокупности, одинаково распределенными нормальными случайными N -векторами с невырожденным распределением $\mathcal{N}_N(\mu, \Sigma)$, тогда критерий отношения правдоподобия для проверки гипотезы независимости H_0 из (6.36) против альтернативы $H_1 = \overline{H_0}$ может быть записан в виде

$$\begin{cases} H_0 & : V \geq V_o(\alpha); \\ H_1 = \overline{H_0} & : V < V_o(\alpha), \end{cases} \quad (6.40)$$

где статистика критерия является V -статистикой [1]:

$$V = \frac{|\hat{\Sigma}|}{\prod_{k=1}^K |\hat{\Sigma}_{kk}|} = \frac{|A|}{\prod_{k=1}^K |A_{kk}|} = \frac{|\mathcal{R}|}{\prod_{k=1}^K |\mathcal{R}_{kk}|} \in [0, 1], \quad (6.41)$$

а порог $V_o = V_o(\alpha)$ при $n \rightarrow +\infty$ определяется по уровню значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ из асимптотического соотношения:

$$V_o(\alpha) = \exp\left(-\frac{F_{\chi_f^2}^{-1}(1-\alpha)}{g}\right), \quad f = \frac{1}{2}\left(N^2 - \sum_{k=1}^K (N^{(k)})^2\right), \quad (6.42)$$

$$g = g(n) = n - \frac{3}{2} - \frac{1}{3} \frac{N^3 - \sum_{k=1}^K (N^{(k)})^3}{N^2 - \sum_{k=1}^K (N^{(k)})^2},$$

$F_{\chi_f^2}^{-1}(\cdot)$ – квантиль χ^2 -распределения с f степенями свободы.

Доказательство. Воспользуемся теоремой 5.1, согласно доказательству которой выборочные среднее \bar{x} и ковариационная матрица $\hat{\Sigma}$ являются оценками максимального правдоподобия и максимизируют функцию правдоподобия $L_N(\mu, \Sigma)$ (логарифмическую функцию правдоподобия $l(\mu, \Sigma) = \ln L_N(\mu, \Sigma)$ – в доказательстве теоремы 5.1) по $(\mu, \Sigma) \in \Theta$:

$$\max_{(\mu, \Sigma) \in \Theta} L_N(\mu, \Sigma) = L_N(\bar{x}, \hat{\Sigma}) = \frac{1}{(2\pi)^{\frac{1}{2}nN} |\hat{\Sigma}|^{\frac{1}{2}n}} e^{-\frac{1}{2}nN}.$$

При истинной гипотезе H_0 получим

$$\begin{aligned} \max_{(\mu, \Sigma) \in \Theta_0} L_N(\mu, \Sigma) &= \prod_{k=1}^K \max_{\mu^{(k)}, \Sigma_{kk}} L_{N^{(k)}}(\mu^{(k)}, \Sigma_{kk}) = \\ &= \prod_{k=1}^K L_{N^{(k)}}(\bar{x}^{(k)}, \hat{\Sigma}_{kk}) = \prod_{k=1}^K \frac{1}{(2\pi)^{\frac{1}{2}nN^{(k)}} |\hat{\Sigma}_{kk}|^{\frac{1}{2}n}} e^{-\frac{1}{2}nN^{(k)}} = \\ &= \frac{1}{(2\pi)^{\frac{1}{2}nN} (\prod_{k=1}^K |\hat{\Sigma}_{kk}|)^{\frac{1}{2}n}} e^{-\frac{1}{2}nN}, \end{aligned}$$

где учтено, что $\bar{x}^{(k)}$ и $\hat{\Sigma}_{kk}$ являются оценками максимального правдоподобия для математического ожидания $\mu^{(k)}$ и ковариационной матрицы Σ_{kk} маргинального нормального распределения $\mathcal{N}_{N^{(k)}}(\mu^{(k)}, \Sigma_{kk})$ и максимизируют соответствующую функцию правдоподобия $L_{N^{(k)}}(\mu^{(k)}, \Sigma_{kk})$, $k = 1, \dots, K$.

Из последних двух соотношений для статистики отношения правдоподобия имеем

$$\lambda = \frac{\max_{(\mu, \Sigma) \in \Theta_0} L_N(\mu, \Sigma)}{\max_{(\mu, \Sigma) \in \Theta} L_N(\mu, \Sigma)} = \left(\frac{|\hat{\Sigma}|}{\prod_{k=1}^K |\hat{\Sigma}_{kk}|} \right)^{\frac{n}{2}} = V^{\frac{n}{2}} -$$

монотонно возрастающая функция от V -статистики

$$V = \frac{|\hat{\Sigma}|}{\prod_{k=1}^K |\hat{\Sigma}_{kk}|},$$

что приводит к эквивалентной записи критерия отношения правдоподобия (6.39), (6.37) в виде (6.40), (6.41).

Покажем, что для V также справедливы представления

$$V = \frac{|A|}{\prod_{k=1}^K |A_{kk}|} = \frac{|\mathcal{R}|}{\prod_{k=1}^K |\mathcal{R}_{kk}|}.$$

Первое из них очевидно и следует из того, что $\hat{\Sigma} = \frac{1}{n}A$. Докажем второе представление. Воспользуемся соотношением, связывающим корреляционную матрицу \mathcal{R} с матрицей $A = (a_{kl})_{k,l=1}^N$:

$$\mathcal{R} = CAC, \quad C = \text{diag} \left\{ \frac{1}{\sqrt{a_{11}}}, \dots, \frac{1}{\sqrt{a_{NN}}} \right\},$$

и разобьем диагональную матрицу C на соответствующие блоки:

$$C = \text{diag}\{C_{11}, \dots, C_{KK}\}, \quad |C| = \prod_{k=1}^K |C_{kk}|,$$

тогда

$$\begin{aligned} \frac{|\mathcal{R}|}{\prod_{k=1}^K |\mathcal{R}_{kk}|} &= \frac{|CAC|}{\prod_{k=1}^K |C_{kk}A_{kk}C_{kk}|} = \frac{|C|^2|A|}{\prod_{k=1}^K |C_{kk}|^2|A_{kk}|} = \\ &= \frac{|A|}{\prod_{k=1}^K |A_{kk}|} = V. \end{aligned}$$

Осталось определить порог критерия $V_o = V_o(\alpha)$. В [1] установлено, что при $n \rightarrow +\infty$ в условиях гипотезы H_0 для распределения вероятностей V -статистики справедливо асимптотическое разложение ($g = g(n) \rightarrow +\infty$):

$$\mathbf{P}_{H_0}\{-g \ln(V) \leq z\} = F_{\chi_f^2}(z) + O(g^{-2}), \quad z \in \mathbb{R},$$

из которого по заданному уровню значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ и определяется порог (6.42). \square

Следствие 6.3. *В условиях теоремы 6.7 критерий отношения правдоподобия для проверки гипотезы о независимости всех N компонент между собой ($K = N$; $N^{(k)} = 1, k = 1, \dots, N$) против общей альтернативы имеет вид*

$$\begin{cases} H_0 & : |\mathcal{R}| \geq V_o(\alpha); \\ H_1 = \overline{H_0} & : |\mathcal{R}| < V_o(\alpha), \end{cases}$$

где $|\mathcal{R}|$ – определитель выборочной корреляционной матрицы, а для порога критерия $V_o = V_o(\alpha)$, $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$, при $n \rightarrow +\infty$ справедливо соотношение

$$V_o(\alpha) = \exp \left(- \left(n - \frac{3}{2} - \frac{N+1}{3} \right)^{-1} F_{\chi_{\frac{N(N-1)}{2}}^2}^{-1} (1 - \alpha) \right).$$

7 Проверка гипотез в многомерном статистическом анализе

7.1 T^2 -статистика Хотеллинга: ее свойства и распределение вероятностей

7.1.1 T^2 -статистика Хотеллинга и ее свойства

Пусть в пространстве \mathbb{R}^N регистрируются наблюдения x_1, \dots, x_n ($x_j \in \mathbb{R}^N$, $j = \overline{1, n}$), являющиеся независимыми одинаково распределенными нормальными случайными векторами

$$\mathcal{L}\{x_j\} = \mathcal{N}_N(\mu, \Sigma), \quad |\Sigma| \neq 0, \quad j = \overline{1, n}.$$

Пусть также зафиксирована некоторая точка $\mu^0 \in \mathbb{R}^N$.

Введем в рассмотрение статистики:

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j,$$
$$S = \frac{1}{n-1} A, \quad \text{где } A = \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T,$$

\bar{x} – выборочное среднее, S – несмещенная оценка ковариационной матрицы ковариации Σ . Везде далее считаем, что $n > N$ и $|S| \neq 0$.

Определение 7.1. T^2 -статистикой Хотеллинга называется следующая статистика:

$$T^2 = T^2(\mu) = n(\bar{x} - \mu)^T S^{-1}(\bar{x} - \mu), \quad (7.1)$$

где $\mu \in \mathbb{R}^N$ – фиксированный N -вектор.

Свойства T^2 -статистики Хотеллига

Свойство 1. $T^2 \geq 0$ и $T^2 = 0 \iff \bar{x} = \mu$.

Доказательство. Обозначим

$$u ::= \sqrt{n}(\bar{x} - \mu).$$

Тогда T^2 -статистика Хотеллинга записывается в виде

$$T^2 = T^2(\mu) = u^T S^{-1} u.$$

Так как $S \succ 0$, $|S| \neq 0$, то $S^{-1} \succ 0$, $|S^{-1}| \neq 0$. Следовательно,

$$T^2 \geq 0 \text{ и } T^2 = 0 \iff \bar{x} = \mu.$$

□

Свойство 2. При $N > 1$ T^2 -статистика пропорциональна оценке квадрата расстояния Махаланобиса.

При $N = 1$ T^2 -статистика пропорциональна квадрату отклонения $(\bar{x} - \mu)^2$.

Доказательство. Функция

$$\rho(\mu, \mu^o) = \sqrt{(\mu - \mu^o)^T \Sigma^{-1} (\mu - \mu^o)},$$

которая носит название расстояние Махаланобиса, определяет расстояние между истинным значением вектора математического ожидания μ и его предполагаемым значением μ^o .

Статистика, вычисленная по случайной выборке $X = \{x_1, \dots, x_n\}$ ($x_j \in \mathbb{R}^N$, $j = \overline{1, n}$),

$$\hat{\rho}(\bar{x}, \mu^o) = \sqrt{(\bar{x} - \mu^o)^T S^{-1} (\bar{x} - \mu^o)} - \quad (7.2)$$

является оценкой расстояния Махаланобиса $\rho(\mu, \mu^o)$.

Сравнивая равенства

$$\hat{\rho}^2(\bar{x}, \mu^o) = (\bar{x} - \mu^o)^T S^{-1} (\bar{x} - \mu^o)$$

и

$$T^2 = T^2(\mu^o) = n(\bar{x} - \mu^o)^T S^{-1} (\bar{x} - \mu^o),$$

видим, что T^2 -статистика пропорциональна оценке квадрата расстояния Махаланобиса.

Для случайной выборки $X = \{x_1, \dots, x_n\}$ из некоторого одномерного распределения вероятностей получаем, полагая в равенстве (7.1) $N = 1$, что T^2 -статистика пропорциональна квадрату отклонения $(\bar{x} - \mu)^2$. \square

Свойство 3. При $\mu^0 = \mathbf{0}_N$ T^2 -статистика инвариантна относительно невырожденного линейного преобразования

$$y_j = Cx_j, \quad C - (N \times N)\text{-матрица, } |C| \neq 0, \quad j = \overline{1, n}.$$

Доказательство. Введем обозначение

$$T_y^2 = T_y^2(\mu^0) \Big|_{\mu^0 = \mathbf{0}_N} = n\bar{y}^T S_y^{-1} \bar{y}, \quad (7.3)$$

где

$$\bar{y} = \frac{1}{n} \sum_{j=1}^n ny_j = \frac{1}{n} \sum_{j=1}^n nCx_j = C\bar{x}, \quad (7.4)$$

$$\begin{aligned} S_y &= \frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y})(y_j - \bar{y})^T = \frac{1}{n-1} \sum_{j=1}^n (Cx_j - C\bar{x})(Cx_j - C\bar{x})^T = \\ &= \frac{1}{n-1} \sum_{j=1}^n C(x_j - \bar{x})(x_j - \bar{x})^T C^T = C \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T C^T = \\ &= CS_x C^T. \end{aligned} \quad (7.5)$$

Подстановкой (7.4) и (7.5) в (7.3) получаем

$$T_y^2 = n\bar{y}^T S_y^{-1} \bar{y} = n\bar{x}^T C^T (C^T)^{-1} S_x^{-1} C^{-1} C\bar{x} = T_x^2.$$

\square

7.1.2 Распределение вероятностей T^2 -статистики Хотеллинга

Определение 7.2. Пусть ξ_1, \dots, ξ_m – независимые в совокупности гауссовские случайные величины с законами распределения вероятностей

$$\mathcal{N}_1(\mu_1, 1), \dots, \mathcal{N}_1(\mu_m, 1).$$

Тогда распределение вероятностей случайной величины

$$\eta = \xi_1^2 + \dots + \xi_m^2$$

называется нецентральным χ^2 -распределением с m степенями свободы и параметром нецентральности τ^2 :

$$\mathcal{L}(\eta) = \chi_{m, \tau^2}^2,$$

где

$$\tau^2 = \mu_1^2 + \dots + \mu_m^2.$$

При $\mu_1 = \dots = \mu_m = 0$ нецентральное χ^2 -распределение называется просто χ^2 -распределением:

$$\chi_{m, \tau^2}^2 \Big|_{\tau^2=0} = \chi_m^2.$$

Определение 7.3. Пусть случайная величина η_1 имеет нецентральное χ^2 -распределение с m степенями свободы и параметром нецентральности τ^2 :

$$\mathcal{L}(\eta_1) = \chi_{m, \tau^2}^2,$$

а случайная величина η_2 не зависит от η_1 и имеет распределение χ_n^2 :

$$\mathcal{L}(\eta_2) = \chi_n^2.$$

тогда распределение вероятностей случайной величины $\eta = \frac{\frac{1}{m}\eta_1}{\frac{1}{n}\eta_2}$ называется нецентральным F -распределением с (m, n) степенями свободы и параметром нецентральности τ^2 :

$$\mathcal{L}(\eta) = F_{m, n, \tau^2}.$$

При $\tau^2 = 0$ имеем центральное F -распределение:

$$\mathcal{L}(\eta) = F_{m, n, 0} =: F_{m, n}.$$

Теорема 7.1 (о распределении T^2 -статистики). Пусть x_1, \dots, x_n ($x_j \in \mathbb{R}^N$, $j = \overline{1, n}$) – независимые в совокупности одинаково распределенные по закону $\mathcal{N}_N(\mu, \Sigma)$, ($|\Sigma| \neq 0$). Пусть $n > N$, а выборочная ковариационная матрица S невырожденная ($|S| \neq 0$). Тогда T^2 -статистика, определенная в (7.1), с условием нормировки имеет следующее нецентральное F -распределение:

$$\mathcal{L} \left\{ \frac{n - N}{N} \cdot \frac{T^2}{n - 1} \right\} = F_{N, n - N, \tau^2}, \quad (7.6)$$

$$\tau^2 = n (\mu - \mu^0)^T \Sigma^{-1} (\mu - \mu^0).$$

Следствие 7.1. При $\mu = \mu^0$ в условиях теоремы $\tau^2 = 0$ и

$$\mathcal{L} \left\{ \frac{n - N}{N} \cdot \frac{T^2}{n - 1} \right\} = F_{N, n - N}. \quad (7.7)$$

7.2 Проверка гипотезы о значении вектора математического ожидания. T^2 -Стьюдента и его оптимальность

Пусть в пространстве \mathbb{R}^N наблюдается случайная выборка $X = \{x_1, \dots, x_n\}$ объема n из невырожденного N -мерного нормального распределения $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$). Истинные значения N -вектора математического ожидания $\mu \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицы Σ неизвестны. Необходимо проверить гипотезы о значении вектора математического ожидания:

$$\begin{aligned} H_0 &: \mu = \mu^0; \\ H_1 = \overline{H_0} &: \mu \neq \mu^0, \end{aligned} \quad (7.8)$$

где $\mu^0 \in \mathbb{R}^N$ – предполагаемое (гипотетическое) значение математического ожидания.

Согласно свойству 2: чем меньше значение статистики $T^2(\mu^0)$, тем «ближе» истинное значение математического ожидания μ (выборочное среднее \bar{x} является оценкой μ) к гипотетическому значению μ^0 . Эти соображения позволяют предложить следующий критерий для проверки гипотез (7.8):

$$\begin{cases} H_0 &: T^2(\mu^0) \leq T_o^2(\alpha), \\ H_1 &: T^2(\mu^0) > T_o^2(\alpha). \end{cases} \quad (7.9)$$

где $T_o^2(\alpha)$ – пороговое значение критерия, подлежащее определению по наперед заданному уровню значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$.

Соотношение (7.7) позволяет найти порог критерия (7.9), вычислив вероятность ошибки первого рода:

$$\begin{aligned} \mathbf{P}\{H_1|H_0\} &= \mathbf{P}\{T^2(\mu^0) > T_o^2|\mu = \mu^0\} = 1 - \mathbf{P}\{T^2(\mu^0) \leq T_o^2|\mu = \mu^0\} = \\ &= 1 - F_{N, n-N} \left(\frac{n-N}{N(n-1)} T_o^2 \right) := \alpha, \end{aligned}$$

откуда

$$T_o^2 = T_o^2(\alpha) = \frac{N(n-1)}{n-N} F_{N, n-N}^{-1}(1-\alpha), \quad (7.10)$$

где $F_{N, n-N}^{-1}(1-\alpha)$ – квантиль уровня $1-\alpha$ от F -распределения Фишера с $(N, n-N)$ степенями свободы.

Следует отметить, что порог $T_o^2 = T_o^2(\alpha)$ позволяет также построить доверительную область уровня $1-\alpha$ для вектора математического ожидания:

$$\mathcal{M}_\alpha = \left\{ \mu \in \mathbb{R}^N : (\bar{x} - \mu)^T S^{-1} (\bar{x} - \mu) \leq \frac{T_o^2(\alpha)}{n} \right\},$$

являющуюся случайным N -мерным эллипсоидом в \mathbb{R}^N с центром в точке \bar{x} , форма и «объем» (мера Лебега) которого определяются матрицей S^{-1} и вероятностью $1-\alpha$, с которой в него «попадает» истинное значение вектора математического ожидания μ :

$$\mathbf{P}\{\mu \in \mathcal{M}_\alpha\} = 1 - \alpha.$$

Исследуем вероятностные свойства критерия (7.9), (7.10), называемого еще T^2 -критерием Стьюдента.

Теорема 7.2. T^2 -критерий (7.9), (7.10), построенный по выборке

$$X = \{x_1, \dots, x_n\}$$

объема $n > N$, образованной независимыми в совокупности, одинаково распределенными N -векторами-наблюдениями с невырожденным нормальным распределением $\mathcal{N}_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$), является критерием отношения правдоподобия для проверки гипотез (7.8) о значении вектора математического ожидания

$$\begin{aligned} H_0 &: \mu = \mu_0, \\ H_1 &: \overline{H_0}, \end{aligned} \quad (7.11)$$

при неизвестной ковариационной матрице Σ , и при любом фиксированном уровне значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ является равномерно наиболее мощным критерием среди всех критериев для проверки гипотез (7.8), мощность которых $w = \mathbf{P}\{H_1|H_1\}$ зависит лишь от величины

$$n(\mu - \mu^o)^T \Sigma^{-1}(\mu - \mu^o)$$

Доказательство. Построим критерий отношения правдоподобия для проверки гипотез (7.8):

$$\begin{aligned} H_0 &: \lambda \geq \lambda_o(\alpha), \\ H_1 &: \lambda < \lambda_o(\alpha). \end{aligned} \quad (7.12)$$

где статистика отношения правдоподобия

$$\lambda = \frac{\max_{(\mu, \Sigma) \in \Theta_0} L_N(\mu, \Sigma)}{\max_{(\mu, \Sigma) \in \Theta} L_N(\mu, \Sigma)}$$

определяется функцией правдоподобия $L_N(\mu, \Sigma)$, максимальное значение которой на множестве всех допустимых значений параметров $\Theta = \{(\mu, \Sigma) : \mu \in R^N, \Sigma = \Sigma^T \succ 0\}$ равно:

$$\max_{(\mu, \Sigma) \in \Theta} L_N(\mu, \Sigma) = L_N(\bar{x}, \hat{\Sigma}) = \frac{1}{(2\pi)^{\frac{1}{2}nN} |\hat{\Sigma}|^{\frac{1}{2}n}} e^{-\frac{1}{2}nN}, \quad \hat{\Sigma} = \frac{1}{n} A.$$

На множестве значений параметров Θ_0 , которому соответствует гипотеза H_0 из (7.8): $\Theta_0 = \{(\mu^o, \Sigma) : \Sigma = \Sigma^T \succ 0\}$, имеем

$$\begin{aligned} \max_{(\mu, \Sigma) \in \Theta_0} L_N(\mu, \Sigma) &= \max_{\Sigma = \Sigma^T \succ 0} L_N(\mu^o, \Sigma) = L_N(\mu^o, \hat{\Sigma}_0) = \\ &= \frac{1}{(2\pi)^{\frac{1}{2}nN} |\hat{\Sigma}_0|^{\frac{1}{2}n}} e^{-\frac{1}{2}nN}, \end{aligned}$$

где

$$\hat{\Sigma}_0 = \frac{1}{n} \sum_{t=1}^n (x_t - \mu^o)(x_t - \mu^o)^T = \frac{1}{n} A + (\bar{x} - \mu^o)(\bar{x} - \mu^o)^T.$$

С учетом полученного соотношения преобразуем статистику отношения правдоподобия:

$$\lambda = \frac{L(\mu^o, \hat{\Sigma}_0)}{L(\bar{x}, \hat{\Sigma})} = \left(\frac{|\hat{\Sigma}|}{|\hat{\Sigma}_0|} \right)^{\frac{n}{2}} = \left(\frac{|A|}{|A + n(\bar{x} - \mu^o)(\bar{x} - \mu^o)^T|} \right)^{\frac{n}{2}}.$$

Далее воспользуемся известными формулами для вычисления определителя блочной матрицы:

$$\begin{vmatrix} E & F \\ G & H \end{vmatrix} = |E| |H - GE^{-1}F|, \quad |E| \neq 0;$$

$$\begin{vmatrix} E & F \\ G & H \end{vmatrix} = |H| |E - FH^{-1}G|, \quad |H| \neq 0.$$

Положим $E ::= 1$, $F ::= \sqrt{n}(\bar{x} - \mu^o)^T$, $G ::= -\sqrt{n}(\bar{x} - \mu^o)$, $H = A$ и получим соотношение

$$\begin{aligned} \left| \begin{array}{cc} 1 & \sqrt{n}(\bar{x} - \mu^o)^T \\ -\sqrt{n}(\bar{x} - \mu^o) & A \end{array} \right| &= |1| |A - (-\sqrt{n}(\bar{x} - \mu^o))1^{-1}\sqrt{n}(\bar{x} - \mu^o)^T| = \\ &= |A + n(\bar{x} - \mu^o)(\bar{x} - \mu^o)^T| = |A| |1 + n(\bar{x} - \mu^o)^T A^{-1}(\bar{x} - \mu^o)|, \end{aligned}$$

с учетом которого статистика отношения правдоподобия

$$\lambda = \left(\frac{1}{1 + \frac{T^2(\mu^o)}{n-1}} \right)^{\frac{n}{2}} -$$

строго убывающая функция от T^2 -статистики $T^2(\mu^o)$ из (7.1), поэтому критерий отношения правдоподобия на ее основе эквивалентен T^2 -критерию Стьюдента (7.9), (7.10).

Доказательство того, что T^2 -критерий (7.9), (7.10) является равномерно наиболее мощным, носит технический характер и основано на анализе распределения вероятностей T^2 -статистики в случае, когда гипотеза H_0 неверна. \square

На практике иногда по выборке $X = \{x_1, \dots, x_n\}$ объема n из невырожденного нормального распределения $N_N(\mu, \Sigma)$ ($|\Sigma| \neq 0$) необходимо проверить еще одну гипотезу о значении вектора математического ожидания — *гипотезу симметрии*, которая состоит в предположении о том, что значения компонент вектора математического ожидания $\mu = (\tilde{\mu}_1, \dots, \tilde{\mu}_N)^T$ равны между собой при альтернативе общего вида:

$$\begin{aligned} H_0 &: \tilde{\mu}_1 = \dots = \tilde{\mu}_N, \\ H_1 &: \exists i \neq j \in \{1, \dots, N\}, \tilde{\mu}_i \neq \tilde{\mu}_j. \end{aligned} \quad (7.13)$$

Очевидно, что гипотеза симметрии H_0 из (7.13) может быть записана в виде $\mu = \tilde{\mu} \mathbf{1}_N$, где $\tilde{\mu} \in \mathbb{R}$ — произвольное общее значение компонент, а $\mathbf{1}_N$ — N -вектор, составленный из единиц.

Сведем гипотезы из (7.13) к гипотезам о значении вектора математического ожидания. Для этого введем в рассмотрение произвольную $((N-1) \times N)$ -матрицу C полного ранга: $\text{rank}(C) = N-1$, строки которой линейно независимы между собой и ортогональны к единичному N -вектору $\mathbf{1}_N$, то есть

$$C \mathbf{1}_N = \mathbf{0}_{N-1}.$$

Осуществим линейное преобразование выборки $X = \{x_1, \dots, x_n\}$ в выборку $Y = \{y_1, \dots, y_n\}$:

$$y_j = Cx_j, \quad j = \overline{1, n}.$$

Наблюдения из выборки Y имеют невырожденное нормальное распределение:

$$\mathcal{L}\{y_j\} = N_{N-1}(C\mu, C\Sigma C^T), \quad |C\Sigma C^T| \neq 0,$$

и при истинной гипотезе H_0 , по построению, их математическое ожидание — нулевой $(N-1)$ -вектор:

$$C\mu = \tilde{\mu} C \mathbf{1}_N = \tilde{\mu} \mathbf{0}_{N-1} = \mathbf{0}_{N-1}.$$

Поэтому проверка гипотез (7.13) эквивалентна проверке гипотезы о том, что математическое ожидание наблюдений из выборки $Y = \{y_1, \dots, y_n\}$ — нулевое, против альтернативы

общего вида. Остается лишь воспользоваться T^2 -критерием (7.1), (7.10), (7.9), положив в нем $N ::= N - 1$; $x_t ::= Cx_t$, $t = \overline{1, n}$; $\mu^o ::= \mathbf{0}_{N-1}$.

Получим следующий критерий для проверки гипотез (7.13):

$$\begin{aligned} H_0 &: n\bar{x}^T C^T (CSC^T)^{-1} C\bar{x} \leq \frac{(N-1)(n-1)}{n-(N-1)} F_{N-1, n-(N-1)}^{-1} (1 - \alpha), \\ H_1 &: n\bar{x}^T C^T (CSC^T)^{-1} C\bar{x} > \frac{(N-1)(n-1)}{n-(N-1)} F_{N-1, n-(N-1)}^{-1} (1 - \alpha). \end{aligned} \quad (7.14)$$

где $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ — наперед заданный уровень значимости. Мощность полученного критерия зависит от матрицы линейного преобразования C : $w = w(C) = \mathbf{P}\{H_1|H_1\}$, и его можно оптимизировать по мощности, выбирая соответствующим образом матрицу C .

7.3 Сравнение векторов математических ожиданий по двум выборкам. Многомерная проблема Беренса – Фишера

7.3.1 Сравнение векторов математических ожиданий при неизвестной одинаковой ковариационной матрице

Пусть в пространстве наблюдений \mathbb{R}^N имеются две случайные выборки

$$X^{(1)} = \{x_1^{(1)}, \dots, x_{N_1}^{(1)}\} \text{ и } X^{(2)} = \{x_1^{(2)}, \dots, x_{N_2}^{(2)}\}$$

соответственно объемов $N_1 > N$ и $N_2 > N$. Наблюдения $\{x_j^{(i)}\}_{j=1}^{N_i}$ из выборки $X^{(i)}$ ($i = 1, 2$) независимы в совокупности и одинаково распределены по нормальному закону $\mathcal{N}_N(\mu_i, \Sigma)$ с математическим ожиданием $\mu_i \in \mathbb{R}^N$ и невырожденной ковариационной $(N \times N)$ -матрицей Σ ($|\Sigma| \neq 0$). Выборка $X^{(1)}$ не зависит от выборки $X^{(2)}$. Необходимо проверить гипотезу о совпадении векторов математических ожиданий μ_1 и μ_2 против альтернативы общего вида:

$$\begin{aligned} H_0 &: \mu_1 = \mu_2, \\ H_1 &: \mu_1 \neq \mu_2. \end{aligned} \quad (7.15)$$

при неизвестном значении одинаковой для обеих выборок ковариационной матрицы Σ .

Для проверки гипотез (7.15) воспользуемся предложенной в предыдущем параграфе T^2 -статистикой, а точнее — леммой. По выборкам $X^{(1)}$ и $X^{(2)}$ определим статистику:

$$y = \frac{1}{\sqrt{\frac{1}{N_1} + \frac{1}{N_2}}} (\bar{x}_{(1)} - \bar{x}_{(2)}), \quad (7.16)$$

а также построим оценку для общей ковариационной матрицы Σ , которая учитывает все $N_1 + N_2$ наблюдений из выборок $X^{(1)}$ и $X^{(2)}$:

$$S = \frac{1}{N_1 + N_2 - 2} ((N_1 - 1)S_{(1)} + (N_2 - 1)S_{(2)}), \quad (7.17)$$

где

$$\bar{x}_{(i)} = \frac{1}{N_i} \sum_{j=1}^{N_i} x_j^{(i)}, \quad S_{(i)} = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (x_j^{(i)} - \bar{x}_{(i)})(x_j^{(i)} - \bar{x}_{(i)})^T -$$

выборочные среднее и ковариационная матрица для выборки $X^{(i)}$ ($i = 1, 2$).

Ранее было доказано, что $\bar{x}_{(i)}$ и $S_{(i)}$ независимы и имеют распределения:

$$\mathcal{L}\{\bar{x}_{(i)}\} = \mathcal{N}_N\left(\mu_i, \frac{1}{N_i}\Sigma\right);$$

$$\mathcal{L}\{(N_i - 1)S_{(i)}\} = W_N(\Sigma, N_i - 1), \quad i = 1, 2,$$

поэтому в силу независимости выборок $X^{(1)}$ и $X^{(2)}$ статистики y и S также независимы, а по свойствам многомерного нормального распределения и распределения Уишарта:

$$\mathcal{L}\{y\} = \mathcal{N}_N \left(\frac{1}{\sqrt{\frac{1}{N_1} + \frac{1}{N_2}}}(\mu_1 - \mu_2), \Sigma \right);$$

$$\mathcal{L}\{(N_1 + N_2 - 2)S\} = W_N(\Sigma, N_1 + N_2 - 2).$$

Если через μ_y обозначить

$$\mu_y = \mathbf{E}\{y\} = \frac{1}{\sqrt{\frac{1}{N_1} + \frac{1}{N_2}}}(\mu_1 - \mu_2),$$

то гипотезы (7.15) примут вид

$$\begin{aligned} H_0 &: \mu_y = \mathbf{0}_N, \\ H_1 &: \mu_y \neq \mathbf{0}_N. \end{aligned} \quad (7.18)$$

Введем в рассмотрение *обобщенную \bar{T}^2 -статистику*:

$$\bar{T}^2 = y^T S^{-1} y, \quad (7.19)$$

основанный на ней критерий для проверки гипотез (7.15) имеет вид

$$\begin{cases} H_0 &: \bar{T}^2 \leq \bar{T}_o^2, \\ H_1 &: \bar{T}^2 > \bar{T}_o^2. \end{cases} \quad (7.20)$$

где порог $\bar{T}_o^2 = \bar{T}_o^2(\alpha)$ определяется по наперед заданному малому значению уровня значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$.

Найдем пороговое значение $\bar{T}_o^2 = \bar{T}_o^2(\alpha)$, применив лемму к \bar{T}^2 -статистике, задаваемой соотношениями (7.19), (7.16), (7.17):

$$\begin{aligned} \mathbf{P}\{H_1|H_0\} &= 1 - \mathbf{P}\{H_0|H_0\} = 1 - \mathbf{P}\{\bar{T}^2 \leq \bar{T}_o^2 | H_0\} = \\ &= 1 - \mathbf{P}\left\{ \frac{N_1 + N_2 - N - 1}{N} \frac{\bar{T}^2}{N_1 + N_2 - 2} \leq \frac{N_1 + N_2 - N - 1}{N} \frac{\bar{T}_o^2}{N_1 + N_2 - 2} \right\} = \\ &= 1 - F_{N, N_1 + N_2 - N - 1} \left(\frac{N_1 + N_2 - N - 1}{N} \frac{\bar{T}_o^2}{N_1 + N_2 - 2} \right) := \alpha, \end{aligned}$$

откуда получаем

$$\bar{T}_o^2 = \bar{T}_o^2(\alpha) = \frac{N(N_1 + N_2 - 2)}{N_1 + N_2 - N - 1} F_{N, N_1 + N_2 - N - 1}^{-1}(1 - \alpha). \quad (7.21)$$

Следует заметить, что критерий (7.20), (7.21) является критерием отношения правдоподобия.

7.3.2 Сравнение векторов математических ожиданий при различных ковариационных матрицах. Многомерная проблема Беренца – Фишера

Обобщим задачу из предыдущего пункта. Пусть в пространстве наблюдений \mathbb{R}^N имеются две случайные выборки

$$X^{(1)} = \{x_1^{(1)}, \dots, x_{N_1}^{(1)}\} \text{ и } X^{(2)} = \{x_1^{(2)}, \dots, x_{N_2}^{(2)}\}$$

соответственно объемов $N_1 > N$ и $N_2 > N$. Наблюдения $\{x_j^{(i)}\}_{j=1}^{N_i}$ из выборки $X^{(i)}$ ($i = 1, 2$) независимы в совокупности и одинаково распределены по нормальному закону $\mathcal{N}_N(\mu_i, \Sigma_i)$ с математическим ожиданием $\mu_i \in \mathbb{R}^N$ и невырожденной ковариационной $(N \times N)$ -матрицей Σ_i ($|\Sigma_i| \neq 0$), $i = 1, 2$. Выборка $X^{(1)}$ не зависит от выборки $X^{(2)}$. Необходимо проверить гипотезы о совпадении векторов математических ожиданий при условии, что ковариационные матрицы Σ_1 и Σ_2 , вообще говоря, различны, а их значения неизвестны:

$$\begin{aligned} H_0 &: \mu_1 = \mu_2, \\ H_1 &: \mu_1 \neq \mu_2. \end{aligned} \tag{7.22}$$

Рассмотрим два случая.

Выборки равного объема ($N_1 = N_2 = n$). Образует новую выборку

$$X = \{x_1, \dots, x_n\}$$

:

$$x_j = x_j^{(1)} - x_j^{(2)}, \quad j = \overline{1, n},$$

где $\{x_j\}_{j=1}^n$ независимы в совокупности и одинаково распределены ($j = \overline{1, n}$):

$$\mathcal{L}\{x_t\} = \mathcal{N}_N(\mu_1 - \mu_2, \Sigma_1 + \Sigma_2).$$

Обозначим: $\mu = \mathbf{E}\{x_t\} = \mu_1 - \mu_2$, тогда гипотезы (7.22) примут эквивалентный вид:

$$\begin{aligned} H_0 &: \mu = \mathbf{0}_N, \\ H_1 &: \mu \neq \mathbf{0}_N. \end{aligned}$$

Построим выборочную ковариационную матрицу по выборке X :

$$S = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T = S_{(1)} + S_{(2)},$$

где

$$\bar{x} = \bar{x}_{(1)} - \bar{x}_{(2)},$$

а $\bar{x}_{(1)}$, $\bar{x}_{(2)}$ и $S_{(1)}$, $S_{(2)}$ — выборочные средние и выборочные ковариационные матрицы из (7.17), вычисленные по выборкам $X^{(1)}$, $X^{(2)}$.

Далее применяем к выборке $X = \{x_j\}_{j=1}^n$ T^2 -критерий Стьюдента, полагая в нем $\mu^o := \mathbf{0}_N$, для T^2 -статистики получим соотношение

$$T^2 = n\bar{x}^T S^{-1} \bar{x} = n(\bar{x}_{(1)} - \bar{x}_{(2)})^T (S_{(1)} + S_{(2)})^{-1} (\bar{x}_{(1)} - \bar{x}_{(2)}),$$

а сам T^2 -критерий имеет вид

$$\begin{cases} H_0 &: T^2 \leq T_o^2(\alpha), \\ H_1 &: T^2 > T_o^2(\alpha). \end{cases}$$

где пороговое значение $T_o^2(\alpha)$ вычисляется по уровню значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ из соотношения:

$$T_o^2(\alpha) = \frac{N(n-1)}{n-N} F_{N, n-N}^{-1}(1-\alpha).$$

Выборки различного объема ($N_1 \neq N_2$; пусть $N_1 < N_2$). В этом случае можно из выборки большего объема отбросить наблюдения, чтобы ее объем стал равен объему меньшей выборки, и воспользоваться критерием, полученным выше для случая равных объемов.

Однако чтобы учесть все наблюдения, целесообразнее поступить следующим образом [1]. Ввести в рассмотрение новую выборку $X = \{x_t\}_{t=1}^{N_1}$:

$$x_t = x_t^{(1)} - \frac{1}{N_2} \sum_{j=1}^{N_2} x_j^{(2)} - \sqrt{\frac{N_1}{N_2}} \left(x_t^{(2)} - \frac{1}{N_1} \sum_{i=1}^{N_1} x_i^{(2)} \right), \quad t = \overline{1, N_1},$$

где $\{x_t\}_{t=1}^{N_1}$ некоррелированы:

$$\begin{aligned} \mathbf{Cov}\{x_t, x_l\} &= \\ &= \mathbf{Cov} \left\{ x_t^{(1)} - \frac{1}{N_2} \sum_{j=1}^{N_2} x_j^{(2)} - \sqrt{\frac{N_1}{N_2}} \left(x_t^{(2)} - \frac{1}{N_1} \sum_{i=1}^{N_1} x_i^{(2)} \right), \right. \\ &\quad \left. x_l^{(1)} - \frac{1}{N_2} \sum_{j=1}^{N_2} x_j^{(2)} - \sqrt{\frac{N_1}{N_2}} \left(x_l^{(2)} - \frac{1}{N_1} \sum_{i=1}^{N_1} x_i^{(2)} \right) \right\} = \\ &= \delta_{tl} \Sigma_1 + \mathbf{Cov} \left\{ \frac{1}{N_2} \sum_{j=1}^{N_2} x_j^{(2)} + \sqrt{\frac{N_1}{N_2}} \left(x_t^{(2)} - \frac{1}{N_1} \sum_{i=1}^{N_1} x_i^{(2)} \right), \right. \\ &\quad \left. \frac{1}{N_2} \sum_{j=1}^{N_2} x_j^{(2)} + \sqrt{\frac{N_1}{N_2}} \left(x_l^{(2)} - \frac{1}{N_1} \sum_{i=1}^{N_1} x_i^{(2)} \right) \right\} = \\ &= \delta_{tl} \Sigma_1 + \mathbf{Cov} \left\{ \sqrt{\frac{N_1}{N_2}} x_t^{(2)} + \left(\frac{1}{N_2} - \frac{1}{\sqrt{N_1 N_2}} \right) \sum_{j=1}^{N_1} x_j^{(2)} + \frac{1}{N_2} \sum_{j=N_1+1}^{N_2} x_j^{(2)}, \right. \\ &\quad \left. \sqrt{\frac{N_1}{N_2}} x_l^{(2)} + \left(\frac{1}{N_2} - \frac{1}{\sqrt{N_1 N_2}} \right) \sum_{j=1}^{N_1} x_j^{(2)} + \frac{1}{N_2} \sum_{j=N_1+1}^{N_2} x_j^{(2)} \right\} = \\ &= \delta_{tl} \Sigma_1 + \delta_{tl} \frac{N_1}{N_2} \Sigma_2 + 2 \left(\frac{1}{N_2} - \frac{1}{\sqrt{N_1 N_2}} \right) \sqrt{\frac{N_1}{N_2}} \Sigma_2 + \\ &\quad + \left(\frac{1}{N_2} - \frac{1}{\sqrt{N_1 N_2}} \right)^2 N_1 \Sigma_2 + \left(\frac{1}{N_2} \right)^2 (N_2 - N_1) \Sigma_2 = \\ &= \delta_{tl} \left(\Sigma_1 + \frac{N_1}{N_2} \Sigma_2 \right), \quad t, l = \overline{1, N_1}, \end{aligned}$$

имеют нормальное распределение:

$$\mathcal{L}\{x_t\} = \mathcal{N}_N \left(\mu_1 - \mu_2, \Sigma_1 + \frac{N_1}{N_2} \Sigma_2 \right), \quad t = \overline{1, N_1},$$

и поэтому независимы в совокупности.

По построенной выборке $X = \{x_t\}_{t=1}^{N_1}$ объема N_1 вычислим T^2 -статистику из $(\mu^o ::= \mathbf{0}_N, n ::= N_1)$:

$$T^2 = N_1 \bar{x}^T S^{-1} \bar{x},$$

где

$$\bar{x} = \frac{1}{N_1} \sum_{t=1}^{N_1} x_t, \quad S = \frac{1}{N_1 - 1} \sum_{t=1}^{N_1} (x_t - \bar{x})(x_t - \bar{x})^T -$$

выборочные среднее и ковариационная матрица по $X = \{x_t\}_{t=1}^{N_1}$.

Далее для проверки гипотез H_0, H_1 применяем к выборке $X = \{x_t\}_{t=1}^{N_1}$ T^2 -критерий Стьюдента:

$$\begin{aligned} H_0 &: T^2 \leq T_o^2(\alpha), \\ H_1 &: T^2 > T_o^2(\alpha). \end{aligned}$$

$$T_o^2(\alpha) = \frac{N(N_1 - 1)}{N_1 - N} F_{N, N_1 - N}^{-1}(1 - \alpha), \quad \alpha = \mathbf{P}\{H_1 | H_0\} \in (0, 1).$$

Замечание 7.3.1. Для построенных в критериев не показано, что они являются оптимальными по мощности. В этом и состоит многомерная проблема Беренса – Фишера. Поэтому если заранее известно, что ковариационные матрицы в выборках совпадают, то лучше использовать критерий из, являющийся критерием отношения правдоподобия.

7.4 Проверка гипотез относительно параметров многомерного нормального распределения

В пространстве \mathbb{R}^N наблюдается случайная выборка $X = \{x_1, \dots, x_n\}$ объема $n > N$ из невырожденного нормального распределения $\mathcal{N}_N(\mu, \Sigma)$ с неизвестными значениями параметров: вектора математического ожидания $\mu \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицы Σ ($|\Sigma| \neq 0$). Рассмотрим всевозможные гипотезы относительно μ и Σ .

7.4.1 Проверка гипотез о значении вектора математического ожидания

Пусть ковариационная матрица Σ ($|\Sigma| \neq 0$) неизвестна. Необходимо проверить гипотезы о том, что математическое ожидание μ совпадает с наперед заданным вектором $\mu^o \in \mathbb{R}^N$:

$$\begin{aligned} H_0 &: \mu = \mu^o; \\ H_1 = \overline{H_0} &: \mu \neq \mu^o, \end{aligned} \tag{7.23}$$

где $\mu^o \in \mathbb{R}^N$ – предполагаемое (гипотетическое) значение математического ожидания.

Данная задача была решена в параграфе 7.2.

7.4.2 Проверка гипотез о значении ковариационной матрицы

Пусть вектор математического ожидания $\mu \in \mathbb{R}^N$ неизвестен. Необходимо проверить гипотезы:

$$\begin{aligned} H_0 &: \Sigma = \Sigma^o; \\ H_1 = \overline{H_0} &: \Sigma \neq \Sigma^o, \end{aligned} \tag{7.24}$$

где Σ^o – наперед заданная невырожденная ковариационная $(N \times N)$ -матрица ($\Sigma^o = (\Sigma^o)^T \succ 0$).

Воспользуемся критерием отношения правдоподобия:

$$\begin{cases} H_0 &: \lambda \geq \lambda_o(\alpha), \\ H_1 &: \lambda < \lambda_o(\alpha). \end{cases}$$

где статистика отношения правдоподобия вычисляется по функции правдоподобия $L_N(\mu, \Sigma)$

$$L_N(\mu, \Sigma) = \prod_{i=1}^n n_N(x_i | \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{1}{2}nN} |\Sigma|^{\frac{1}{2}n}} \exp \left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right); \tag{7.25}$$

аналогично доказательству теоремы 7.2:

$$\begin{aligned}\lambda &= \frac{\max_{\mu} L_N(\mu, \Sigma^o)}{\max_{\mu, \Sigma} L_N(\mu, \Sigma)} = \frac{L_N(\bar{x}, \Sigma^o)}{L_N(\bar{x}, \hat{\Sigma})} = \\ &= \frac{(2\pi)^{-\frac{nN}{2}} |\Sigma^o|^{-n/2} \exp\left(-\frac{1}{2} \sum_{j=1}^n (x_j - \bar{x})^T (\Sigma^o)^{-1} (x_j - \bar{x})\right)}{(2\pi)^{-\frac{nN}{2}} |\hat{\Sigma}|^{-n/2} e^{-\frac{nN}{2}}} = \\ &= \left(\frac{e}{n}\right)^{\frac{nN}{2}} |\Sigma^o|^{-n/2} |A|^{n/2} \exp\left(-\frac{1}{2} \text{tr}((\Sigma^o)^{-1} A)\right),\end{aligned}$$

где \bar{x} – выборочное среднее:

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j,$$

$\hat{\Sigma}$ – выборочная ковариационная матрица:

$$\hat{\Sigma} = \frac{1}{n} A, \quad A = \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T.$$

Построим критерий, эквивалентный критерию отношения правдоподобия:

$$\begin{cases} H_0 & : -2 \ln \lambda \leq \Delta, \\ H_1 & : -2 \ln \lambda > \Delta. \end{cases}$$

где относительно статистики $-2 \ln \lambda$ установлено [1]:

$$\mathcal{L}_{H_0}\{-2 \ln \lambda\} \rightarrow \chi_f^2, \quad n \rightarrow +\infty,$$

где χ_f^2 – χ^2 -распределение с $f = \frac{N(N+1)}{2}$ степенями свободы. Этот факт позволяет при «большом» объеме выборки n ($n \rightarrow +\infty$) по наперед заданному уровню значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ найти пороговое значение $\Delta = \Delta(\alpha)$:

$$\Delta(\alpha) = F_{\chi_f^2}^{-1}(1 - \alpha), \quad f = \frac{N(N+1)}{2}.$$

7.4.3 Проверка гипотез о совпадении многомерного нормального распределения с наперед заданным многомерным нормальным распределением

Необходимо проверить гипотезы:

$$\begin{aligned} H_0 & : \mu = \mu^o, \quad \Sigma = \Sigma^o, \\ H_1 & : \mu \neq \mu^o \text{ и/или } \Sigma \neq \Sigma^o. \end{aligned}$$

где μ^o и Σ^o – наперед заданные гипотетические значения N -вектора математического ожидания ($\mu^o \in \mathbb{R}^N$) и ковариационной ($N \times N$)-матрицы ($\Sigma^o = (\Sigma^o)^T \succ 0$).

Критерий отношения правдоподобия имеет вид

$$\begin{cases} H_0 & : \lambda \geq \lambda_o(\alpha), \\ H_1 & : \lambda < \lambda_o(\alpha). \end{cases}$$

где

$$\lambda = \frac{L_N(\mu^o, \Sigma^o)}{\max_{\mu, \Sigma} L_N(\mu, \Sigma)} = \frac{L_N(\mu^o, \Sigma^o)}{L_N(\bar{x}, \hat{\Sigma})} =$$

$$\begin{aligned}
&= \frac{(2\pi)^{-\frac{nN}{2}} |\Sigma^o|^{-n/2} \exp\left(-\frac{1}{2} \sum_{t=1}^n (x_t - \mu^o)^T (\Sigma^o)^{-1} (x_t - \mu^o)\right)}{(2\pi)^{-\frac{nN}{2}} \left|\frac{1}{n} A\right|^{-n/2} e^{-\frac{nN}{2}}} = \\
&= \left(\frac{e}{n}\right)^{\frac{nN}{2}} |\Sigma^o|^{-n/2} |A|^{n/2} \exp\left(-\frac{1}{2} (x_t - \mu^o)^T (\Sigma^o)^{-1} (x_t - \mu^o)\right) = \\
&= \left(\frac{e}{n}\right)^{\frac{nN}{2}} |\Sigma^o|^{-n/2} |A|^{n/2} \exp\left(-\frac{1}{2} \text{tr}((\Sigma^o)^{-1} A) - \right. \\
&\quad \left. -\frac{n}{2} (\bar{x} - \mu^o)^T (\Sigma^o)^{-1} (\bar{x} - \mu^o)\right).
\end{aligned}$$

Используем эквивалентный критерий:

$$\begin{cases} H_0 & : -2 \ln \lambda \leq \Delta, \\ H_1 & : -2 \ln \lambda > \Delta. \end{cases}$$

где для статистики $-2 \ln \lambda$ также установлено предельное распределение [1]:

$$\mathcal{L}_{H_0}\{-2 \ln \lambda\} \rightarrow \chi_f^2, \quad f = \frac{N(N+1)}{2} + N, \quad n \rightarrow +\infty,$$

что позволяет определить при $n \rightarrow +\infty$ пороговое значение $\Delta = \Delta(\alpha)$, $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$:

$$\Delta(\alpha) = F_{\chi_f^2}^{-1}(1 - \alpha), \quad f = \frac{N(N+1)}{2} + N.$$

7.5 Проверка гипотез относительно нескольких выборок из многомерных нормальных распределений

Пусть в пространстве \mathbb{R}^N наблюдается $L \geq 2$ независимых в совокупности случайных выборок $X^{(i)} = \{x_1^{(i)}, \dots, x_{n_i}^{(i)}\}$, $i = \overline{1, L}$. Выборка $X^{(i)}$ объема $n_i > N$ ($i = \overline{1, L}$) образована независимыми в совокупности, одинаково распределенными N -векторами-наблюдениями, имеющими невырожденное нормальное распределение $\mathcal{N}_N(\mu_i, \Sigma_i)$ с математическим ожиданием $\mu_i \in \mathbb{R}^N$ и ковариационной $(N \times N)$ -матрицей Σ_i ($|\Sigma_i| \neq 0$).

7.5.1 Гипотеза о совпадении векторов математических ожиданий при неизвестной одинаковой ковариационной матрице

Проверим гипотезу о равенстве между собой векторов математических ожиданий μ_1, \dots, μ_L против альтернативы общего вида:

$$\begin{aligned} H_0 & : \mu_1 = \dots = \mu_L, \\ H_1 & : \overline{H_0}. \end{aligned}$$

при неизвестном значении одинаковой для всех выборок ковариационной матрицы:

$$\Sigma_1 = \dots = \Sigma_L = \Sigma \quad (|\Sigma| \neq 0).$$

Запишем функцию правдоподобия для объединенной выборки

$$X = \bigcup_{i=1}^L X^{(i)}$$

объема

$$n = \sum_{i=1}^L n_i :$$

$$L(\mu_1, \Sigma_1; \dots; \mu_L, \Sigma_L) = \prod_{i=1}^L L_N^{(i)}(\mu_i, \Sigma_i) =$$

$$= (2\pi)^{-\frac{nN}{2}} \prod_{i=1}^L |\Sigma_i|^{-n_i/2} \cdot \exp \left(-\frac{1}{2} \sum_{i=1}^L \sum_{j=1}^{n_i} (x_j^{(i)} - \mu_i)^T \Sigma_i^{-1} (x_j^{(i)} - \mu_i) \right),$$

где учтено, что выборки $X^{(1)}, \dots, X^{(L)}$ независимы между собой, а функция правдоподобия для i -й выборки имеет вид

$$L_N^{(i)}(\mu_i, \Sigma_i) = \prod_{j=1}^{n_i} n_N(x_j^{(i)} | \mu_i, \Sigma_i) =$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}n_i N} |\Sigma_i|^{\frac{1}{2}n_i}} \exp \left(-\frac{1}{2} \sum_{j=1}^{n_i} (x_j^{(i)} - \mu_i)^T \Sigma_i^{-1} (x_j^{(i)} - \mu_i) \right).$$

Воспользуемся критерием отношения правдоподобия и определим статистику отношения правдоподобия:

$$\lambda = \frac{\max_{\mu, \Sigma} L(\mu_1, \Sigma_1; \dots; \mu_L, \Sigma_L)}{\max_{\mu_1, \dots, \mu_L, \Sigma} L(\mu_1, \Sigma; \dots; \mu_L, \Sigma)} =$$

$$= \frac{(2\pi)^{-\frac{nN}{2}} |\frac{1}{n}A|^{-n/2} e^{-\frac{nN}{2}}}{(2\pi)^{-\frac{nN}{2}} |\frac{1}{n} \sum_{i=1}^L A_i|^{-n/2} e^{-\frac{nN}{2}}} = \left(\frac{|\sum_{i=1}^L A_i|}{|A|} \right)^{\frac{n}{2}},$$

где

$$\bar{x} = \frac{1}{n} \sum_{i=1}^L \sum_{j=1}^{n_i} x_j^{(i)} -$$

выборочное среднее,

$$\hat{\Sigma}_0 = \frac{1}{n} A, \quad A = \sum_{i=1}^L \sum_{j=1}^{n_i} (x_j^{(i)} - \bar{x})(x_j^{(i)} - \bar{x})^T -$$

выборочная ковариационная матрица, вычисленные по объединенной выборке X объема n и являющиеся при истинной гипотезе H_0 оценками максимального правдоподобия для общих вектора математического ожидания μ и ковариационной матрицы Σ ;

$$\bar{x}_{(i)} = \frac{1}{n_i} \sum_{j=1}^{n_i} x_j^{(i)}, \quad i = \overline{1, L}, -$$

выборочные средние,

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n A_i, \quad A_i = \sum_{j=1}^{n_i} (x_j^{(i)} - \bar{x}_{(i)})(x_j^{(i)} - \bar{x}_{(i)})^{\mathbf{T}}, \quad i = \overline{1, L}, -$$

выборочная ковариационная матрица, являющиеся оценками максимального правдоподобия для математических ожиданий μ_1, \dots, μ_L и общей для всех выборок $X^{(1)}, \dots, X^{(L)}$ ковариационной матрицы Σ в общем случае (когда гипотеза H_0 , вообще говоря, неверна).

Вместо статистики λ воспользуемся так называемой U -статистикой:

$$U = \lambda_{\frac{2}{n}} = \frac{|\sum_{i=1}^L A_i|}{|A|},$$

для которой известно [1]:

$$\mathcal{L}_{H_0}\{-g \ln U\} \rightarrow \chi_f^2, \quad n_i \rightarrow \infty, \quad i = \overline{1, L} \quad (n \rightarrow +\infty),$$

где $g = n - L - \frac{N-L}{2} - 1$, $f = N(L-1)$. Тогда критерий отношения правдоподобия

$$\begin{cases} H_0 & : \lambda \geq \lambda_o(\alpha), \\ H_1 & : \lambda < \lambda_o(\alpha). \end{cases}$$

запишется в эквивалентном виде:

$$\begin{cases} H_0 & : -g \ln U \leq \Delta(\alpha), \\ H_1 & : -g \ln U > \Delta(\alpha). \end{cases}$$

где порог критерия $\Delta = \Delta(\alpha)$ определяется по наперед заданному уровню значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$:

$$\Delta(\alpha) = F_{\chi_f^2}^{-1}(1 - \alpha), \quad f = N(L-1).$$

В случае двух выборок ($L = 2$), воспользовавшись соотношениями

$$A_1 + A_2 = A + n_1(\bar{x} - \bar{x}_{(1)})(\bar{x} - \bar{x}_{(1)})^{\mathbf{T}} + n_2(\bar{x} - \bar{x}_{(2)})(\bar{x} - \bar{x}_{(2)})^{\mathbf{T}};$$

$$n\bar{x} = n_1\bar{x}_{(1)} + n_2\bar{x}_{(2)},$$

и формулой блочного определителя, по аналогии с доказательством теоремы 7.2 для U -статистики получим

$$\begin{aligned} U &= \frac{|A_1 + A_2|}{|A|} = \\ &= \frac{|A_1 + A_2|}{|A_1 + A_2 - n_1(\bar{x} - \bar{x}_{(1)})(\bar{x} - \bar{x}_{(1)})^{\mathbf{T}} - n_2(\bar{x} - \bar{x}_{(2)})(\bar{x} - \bar{x}_{(2)})^{\mathbf{T}}|} = \\ &= \frac{|A_1 + A_2|}{|A_1 + A_2 - \frac{n_1 n_2}{n_1 + n_2}(\bar{x}_{(1)} - \bar{x}_{(2)})(\bar{x}_{(1)} - \bar{x}_{(2)})^{\mathbf{T}}|} = \\ &= \frac{|A_1 + A_2|}{|A_1 + A_2| |1 + \frac{n_1 n_2}{n_1 + n_2}(\bar{x}_{(1)} - \bar{x}_{(2)})^{\mathbf{T}}(A_1 + A_2)^{-1}(\bar{x}_{(1)} - \bar{x}_{(2)})|} = \\ &= \frac{1}{1 + \frac{\bar{T}^2}{n_1 + n_2 - 2}} - \end{aligned}$$

выражается через обобщенную \bar{T}^2 -статистику из (7.19), и построенный ранее в п. 7.3.1 на ее основе критерий (7.20), (7.21) также является критерием отношения правдоподобия. Более того, он обладает очевидным преимуществом перед критерием, использующем U -статистику: пороговое значение (7.21) в нем не требует асимптотики по объемам выборок. Поэтому для проверки гипотезы о совпадении векторов математических ожиданий в случае двух выборок целесообразнее использовать критерий из п. 7.3.1.

7.5.2 Проверка гипотезы о равенстве ковариационных матриц

Проверим гипотезы о равенстве ковариационных матриц:

$$\begin{aligned} H_0 &: \Sigma_1 = \dots = \Sigma_L, \\ H_1 &: \overline{H_0}. \end{aligned}$$

Статистика отношения правдоподобия в этом случае имеет вид

$$\begin{aligned} \lambda &= \frac{\max_{\mu_1, \dots, \mu_L, \Sigma} L(\{\mu_i, \Sigma_i = \Sigma\}_{i=1}^L)}{\max_{\{\mu_i, \Sigma_i\}_{i=1}^L} L(\{\mu_i, \Sigma_i\}_{i=1}^L)} = \\ &= \frac{\max_{\mu_1, \dots, \mu_L, \Sigma} L(\{\mu_i, \Sigma_i = \Sigma\}_{i=1}^L)}{\prod_{i=1}^L \max_{\mu_i, \Sigma_i} L_N^{(i)}(\mu_i, \Sigma_i)} = \frac{(2\pi)^{-\frac{nN}{2}} \left| \frac{1}{n} \sum_{i=1}^L A_i \right|^{-n/2} e^{-\frac{nN}{2}}}{\prod_{i=1}^L (2\pi)^{-\frac{n_i N}{2}} \left| \frac{1}{n_i} A_i \right|^{-n_i/2} e^{-\frac{n_i N}{2}}} = \\ &= \frac{n^{\frac{nN}{2}}}{\prod_{i=1}^L n_i^{\frac{n_i N}{2}}} \frac{\prod_{i=1}^L |A_i|^{n_i/2}}{\left| \sum_{i=1}^L A_i \right|^{n/2}}. \end{aligned}$$

Введем в рассмотрение «подправленную» статистику, предложенную Бартлеттом [1]:

$$\lambda^* = \frac{m^{\frac{mN}{2}}}{\prod_{i=1}^L m_i^{\frac{m_i N}{2}}} \frac{\prod_{i=1}^L |A_i|^{m_i/2}}{\left| \sum_{i=1}^L A_i \right|^{m/2}},$$

где $m_i = n_i - 1$ ($i = \overline{1, L}$), $m = \sum_{i=1}^L m_i = n - L$, и построим критерий, асимптотически эквивалентный критерию отношения правдоподобия:

$$\begin{cases} H_0 &: -2\rho \ln \lambda^* \leq \Delta(\alpha), \\ H_1 &: -2\rho \ln \lambda^* > \Delta(\alpha). \end{cases}$$

где

$$\rho = 1 - \left(\sum_{i=1}^L \frac{1}{m_i} - \frac{1}{m} \right) \frac{2N^2 + 3N - 1}{6(N+1)(L-1)}$$

и установлено [1]:

$$\mathcal{L}_{H_0}\{-2\rho \ln \lambda^*\} \rightarrow \chi_f^2, \quad f = \frac{L-1}{2}N(N+1), \quad n_i \rightarrow \infty, \quad i = \overline{1, L},$$

что позволяет по уровню значимости $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ определить порог:

$$\Delta(\alpha) = F_{\chi_f^2}^{-1}(1 - \alpha), \quad f = \frac{L-1}{2}N(N+1).$$

7.5.3 Гипотеза об эквивалентности нормальных распределений (гипотеза однородности)

Проверим гипотезы о том, что наблюдения из всех выборок имеют одно и то же нормальное распределение:

$$\begin{aligned} H_0 &: \mu_1 = \dots = \mu_L, \quad \Sigma_1 = \dots = \Sigma_L, \\ H_1 &: \overline{H_0}. \end{aligned}$$

Запишем статистику отношения правдоподобия:

$$\begin{aligned}\lambda &= \frac{\max_{\mu, \Sigma} L(\{\mu_i = \mu, \Sigma_i = \Sigma\}_{i=1}^L)}{\max_{\{\mu_i, \Sigma_i\}_{i=1}^L} L(\{\mu_i, \Sigma_i\}_{i=1}^L)} = \\ &= \frac{(2\pi)^{-\frac{nN}{2}} |\frac{1}{n}A|^{-n/2} e^{-\frac{nN}{2}}}{\prod_{i=1}^L (2\pi)^{-\frac{n_i N}{2}} |\frac{1}{n_i}A_i|^{-n_i/2} e^{-\frac{n_i N}{2}}} = \frac{n^{\frac{nN}{2}}}{\prod_{i=1}^L n_i^{\frac{n_i N}{2}}} \frac{\prod_{i=1}^L |A_i|^{n_i/2}}{|A|^{n/2}}.\end{aligned}$$

Определим соответствующую «подправленную» статистику:

$$\bar{\lambda}^* = \frac{m^{\frac{mN}{2}}}{\prod_{i=1}^L m_i^{\frac{m_i N}{2}}} \frac{\prod_{i=1}^L |A_i|^{m_i/2}}{|A|^{m/2}},$$

и построим критерий, асимптотически эквивалентный критерию отношения правдоподобия и имеющий вид

$$\begin{cases} H_0 & : -2\rho^* \ln \bar{\lambda}^* \leq \Delta(\alpha), \\ H_1 & : -2\rho^* \ln \bar{\lambda}^* > \Delta(\alpha). \end{cases}$$

где

$$\rho^* = 1 - \left(\sum_{i=1}^L \frac{1}{m_i} - \frac{1}{m} \right) \frac{2N^2 + 3N - 1}{6(N+3)(L-1)} - \frac{1}{m} \frac{N-L+2}{N+3}$$

и также установлено [1]:

$$\mathcal{L}_{H_0}\{-2\rho^* \ln \bar{\lambda}^*\} \rightarrow \chi_f^2, \quad f = \frac{L-1}{2}N(N+3), \quad n_i \rightarrow \infty, \quad i = \overline{1, L},$$

что позволяет по $\alpha = \mathbf{P}\{H_1|H_0\} \in (0, 1)$ определить порог:

$$\Delta(\alpha) = F_{\chi_f^2}^{-1}(1 - \alpha), \quad f = \frac{L-1}{2}N(N+3).$$

Список литературы

- [1] *Андерсон, Т.* Введение в многомерный статистический анализ /Т. Андерсон. М. : Физматгиз, 1963.
- [2] *Андерсон, Т.* Статистический анализ временных рядов /Т. Андерсон. М. : Мир, 1976.
- [3] *Бокс, Дж.* Анализ временных рядов. Прогноз и управление /Дж. Бокс, Г. Дженкинс. М. : Мир, 1974. Вып. 1.
- [4] *Боровков, А.А.* Теория вероятностей /А. А. Боровков. М. : Наука, 1998.
- [5] *Гантмахер, Ф.Р.* Теория матриц /Ф.Р. Гантмахер. М. : Наука, 1967.
- [6] *Гнеденко, Б.В.* Курс теории вероятностей /Б. В. Гнеденко М.: Физматгиз, 1988.
- [7] *Ивченко Г. И.* Введение с математическую статистику / Г. И. Ивченко, Ю. И. Медведов М.: Издательство ЛКИ, 2014.
- [8] *Харин, Ю. С.* Теория вероятностей, математическая и прикладная статистика / Ю. С. Харин, Н. М. Зуев, Е. Е. Жук. — Минск : БГУ, 2011.
- [9] *Хорн, Р.* Матричный анализ / Р. Хорн, Ч. Джонсон. М. : Мир, 1989.
- [10] *Ширяев, А.Н.* Вероятность /А.Н. Ширяев. М. : Наука, 1989.
- [11] *Харин, Ю.С.* Теория вероятностей /Ю.С. Харин, Н.М. Зуев. Мн. : БГУ, 2004.
- [12] *Харин, Ю. С.* Теория вероятностей, математическая и прикладная статистика / Ю. С. Харин, Н. М. Зуев, Е. Е. Жук. – Минск : БГУ, 2011.