

**ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«ОРЕНБУРГСКИЙ ГОСУДАРСТВЕННЫЙ АГРАРНЫЙ УНИВЕРСИТЕТ»**

**МЕТОДИЧЕСКИЕ УКАЗАНИЯ ДЛЯ ОБУЧАЮЩИХСЯ
ПО ОСВОЕНИЮ ДИСЦИПЛИНЫ**

Б1.В.ДВ.12.01 ОСНОВЫ НАУЧНЫХ ИССЛЕДОВАНИЙ

Направление подготовки (специальность):

27.03.04 Управление в технических системах

Профиль образовательной программы:

Интеллектуальные системы обработки информации и управления

Форма обучения: заочная

СОДЕРЖАНИЕ

1. Конспект лекций	3
1.1 Лекция № 1 Оптимизационные задачи. Основные методы их решения. Марковские процессы, их приложения к решению инженерных задач	3
1.2 Лекция № 2 Теоретические основы обработки экспериментальных данных. Корреляционно-регрессионный анализ	18
2. Методические указания по проведению практических занятий	41
2.1 Практическое занятие № ПЗ-1 Методологическая основа научно-исследовательской работы	41
2.2 Практическое занятие № ПЗ-2 Математическое моделирование в инженерных исследованиях. Основные понятия и методы математической обработки экспериментальных данных	42
2.3 Практическое занятие № ПЗ-3 Основы корреляционно-регрессионного анализа	46

1. КОНСПЕКТ ЛЕКЦИЙ

1. 1 Лекция №1 (2 часа).

Тема: «Оптимизационные задачи. Основные методы их решения. Марковские процессы, их приложения к решению инженерных задач»

1.1.1 Вопросы лекции:

1. Оптимизация как цель математического моделирования. Виды оптимизационных задач
2. Задача линейного программирования (ЗЛП). Графический метод решения ЗЛП
3. Симплекс-метод решения ЗЛП
4. Простейший поток, его свойства. Классификация потоков. Марковские цепи, их свойства
5. Марковские процессы в инженерной практике

1.1.2 Краткое содержание вопросов:

1. Оптимизация как цель математического моделирования. Виды оптимизационных задач

Многие задачи, с которыми приходится иметь дело в повседневной практике, являются многовариантными. Среди множества вариантов (в условиях рыночных отношений) приходится отыскивать наилучшие, при ограничениях, налагаемых на природные, экономические и технологические возможности. В связи с этим возникла необходимость применять для анализа и синтеза экономических ситуаций и систем математические методы и современную вычислительную технику? Такие методы объединяются под общим названием — математическое программирование.

Основные понятия

Математическое программирование — область математики, разрабатывающая теорию и численные методы решения многомерных экстремальных задач с ограничениями, т. е. задач на экстремум функции многих переменных с ограничениями на область изменения этих переменных.

Функцию, экстремальное значение которой нужно найти в условиях экономических возможностей, называют *целевой*, *показателем эффективности* или *критерием оптимальности*. Экономические возможности формализуются в виде *системы ограничений*. Все это составляет математическую модель. *Математическая модель задачи* — это отражение оригинала в виде функций, уравнений, неравенств, цифр и т. д. Модель задачи математического программирования включает:

1) совокупность неизвестных величин, действуя на которые, систему можно совершенствовать. Их называют *планом задачи* (вектором управления, решением, управлением, стратегией, поведением и др.);

2) целевую функцию (функцию цели, показатель эффективности, критерий оптимальности, функционал задачи и др.). Целевая функция позволяет выбирать наилучший вариант - из множества возможных. Наилучший вариант доставляет целевой функции экстремальное значение. Это может быть прибыль, объем выпуска или реализации, затраты производства, издержки обращения, уровень обслуживания или дефицитности, число комплектов, отходы и т. д.;

Эти условия следуют из ограниченности ресурсов, которыми располагает общество в любой момент времени, из необходимости удовлетворения насущных потребностей, из условий производственных и технологических процессов. Ограниченными являются не

только материальные, финансовые и трудовые ресурсы. Таковыми могут быть возможности технического, технологического и вообще научного потенциала. Нередко потребности превышают возможности их удовлетворения. Математические ограничения выражаются в виде уравнений и неравенств. Их совокупность образует *область допустимых решений (область экономических возможностей)*. План, удовлетворяющий системе ограничений задачи, называется *допустимым*. Допустимый план, доставляющий функции цели экстремальное значение, называется *оптимальным*. Оптимальное решение, вообще говоря, не обязательно единственно, возможны случаи, когда оно не существует, имеется конечное или бесчисленное множество оптимальных решений.

Один из разделов математического программирования - *линейным программированием*. Методы и модели линейного программирования широко применяются при оптимизации процессов во всех отраслях народного хозяйства: при разработке производственной программы предприятия, распределении ее по исполнителям, при размещении заказов между исполнителями и по временным интервалам, при определении наилучшего ассортимента выпускаемой продукции, в задачах перспективного, текущего и оперативного планирования и управления в задачах развития и размещения производительных сил, баз и складов систем обращения материальных ресурсов и т. д. Особенно широкое применение методы и модели линейного программирования получили при решении задач экономии ресурсов (выбор ресурсосберегающих технологий, составление смесей, раскрой материалов), производственно-транспортных и других задач.

Начало линейному программированию было положено в 1939 г. советским математиком-экономистом Л. В. Канторовичем в работе «Математические методы организации и планирования производства». Появление этой работы открыло новый этап в применении математики в экономике. Спустя десять лет американский математик Дж. Данциг разработал эффективный метод решения данного класса задач — симплекс-метод. Общая идея *симплексного метода (метода последовательного улучшения плана)* для решения ЗЛП состоит в следующем:

- 1) умение находить начальный опорный план;
- 2) наличие признака оптимальности опорного плана;
- 3) умение переходить к нехудшему опорному плану.

2. Задача линейного программирования (ЗЛП). Графический метод решения ЗЛП

Постановка задачи линейного программирования и свойства ее решений

Линейное программирование — раздел математического программирования, применяемый при разработке методов отыскания экстремума линейных функций нескольких переменных при линейных дополнительных ограничениях, налагаемых на переменные. По типу решаемых задач его методы разделяются на универсальные и специальные. С помощью универсальных методов могут решаться любые *задачи линейного программирования (ЗЛП)*. Специальные методы учитывают особенности модели задачи, ее целевой функции и системы ограничений.

Особенностью задач линейного программирования является то, что экстремума целевая функция достигает на границе области допустимых решений. Классические же методы дифференциального исчисления связаны с нахождением экстремумов функции во внутренней точке области допустимых значений. Отсюда — необходимость разработки новых методов.

Формы записи задачи линейного программирования:

Общей задачей линейного программирования называют задачу

$$\max(\min) Z = \sum_{j=1}^n c_j x_j \quad (2.1)$$

при ограничениях

$$\sum_{j=1}^n a_{ij}x_j \leq b_i \quad (i = 1, \dots, m_1) \quad (2.2)$$

$$\sum_{j=1}^n a_{ij}x_j = b_i \quad (i = m_1 + 1, \dots, m_2) \quad (2.3)$$

$$\sum_{j=1}^n a_{ij}x_j \geq b_i \quad (i = m_2 + 1, \dots, m) \quad (2.4)$$

$$x_j \geq 0 \quad (j = \overline{1, n_1}) \quad (2.5)$$

$$x_j \text{ - произвольные} \quad (j = n_1 + 1, \dots, n) \quad (2.6)$$

где c_j, a_{ij}, b_i - заданные действительные числа; (2.1) – целевая функция; (2.1) – (2.6) – ограничения; $\vec{x} = (x_1, \dots, x_n)$ - план задачи.

Пусть ЗЛП представлена в следующей записи:

$$\max Z = cx \quad (2.7)$$

$$A_1x_1 + A_2x_2 + \dots + A_nx_n = A_0 \quad (2.8)$$

$$x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0 \quad (2.9)$$

Чтобы задача (2.7) – (2.8) имела решение, системе её ограничений (2.8) должна быть совместной. Это возможно, если r этой системы не больше числа неизвестных n . Случай $r > n$ вообще невозможен. При $r = n$ система имеет единственное решение, которое будет при $x_j \geq 0 \quad (j = \overline{1, n})$ оптимальным. В этом случае проблема выбора оптимального решения теряет смысл. Выясним структуру координат угловой точки многогранных решений.

Пусть $r < n$. В этом случае система векторов A_1, A_2, \dots, A_n содержит базис — максимальную линейно независимую подсистему векторов, через которую любой вектор системы может быть выражен как ее линейная комбинация. Базисов, вообще говоря, может быть несколько, но не более C_n^r . Каждый из них состоит точно из r векторов. Переменные ЗЛП, соответствующие r векторам базиса, называют, как известно, *базисными* и обозначают БП. Остальные $n - r$ переменных будут *свободными*, их обозначают СП. Не ограничивая общности, будем считать, что базис составляют первые m векторов A_1, A_2, \dots, A_m . Этому базису соответствуют базисные переменные x_1, x_2, \dots, x_m , а свободными будут переменные $x_{m+1}, x_{m+2}, \dots, x_n$.

Если свободные переменные приравнять нулю, а базисные переменные при этом примут неотрицательные значения, то полученное частное решение системы (8) называют *опорным решением (планом)*.

Теорема. Если система векторов A_1, A_2, \dots, A_n содержит m линейно независимых векторов A_1, A_2, \dots, A_m , то допустимый план $\vec{x} = (x_1, x_2, \dots, x_m, \underbrace{0; 0; \dots; 0}_{n-m})$ является крайней точкой многогранника планов. (2.10)

Теорема. Если ЗЛП имеет решение, то целевая функция достигает экстремального значения хотя бы в одной из крайних точек многогранника решений. Если же целевая функция достигает экстремального значения более чем в одной крайней точке, то она достигает того же значения в любой точке, являющейся их выпуклой линейной комбинацией.

Графический способ решения ЗЛП

Геометрическая интерпретация экономических задач дает возможность наглядно представить, их структуру, выявить особенности и открывает пути исследования более сложных свойств. ЗЛП с двумя переменными всегда можно решить графически. Однако

Вектор $\vec{c} = (c_1, c_2)$ перпендикулярен к прямым $Z = \text{const}$ семейства $c_1x_1 + c_2x_2 = Z$.

Из геометрической интерпретации элементов ЗЛП вытекает следующий порядок ее графического решения.

1. С учетом системы ограничений строим область допустимых решений Ω .
2. Строим вектор $\vec{c} = (c_1, c_2)$ наискорейшего возрастания целевой функции — вектор градиентного направления.

3. Проводим произвольную линию уровня $Z = Z_0$.

4. При решении задачи на максимум перемещаем линию уровня $Z = Z_0$ в направлении вектора \vec{c} так, чтобы она касалась области допустимых решений в ее крайнем положении (крайней точке). В случае решения задачи на минимум линию уровня $Z = Z_0$ перемещают в антиградиентном направлении.

5. Определяем оптимальный план $\vec{x}^* = (x_1^*, x_2^*)$ и экстремальное значение целевой функции $Z^* = z(\vec{x}^*)$.

3. Симплексный метод решение ЗЛП

Общая идея симплексного метода (метода последовательного улучшения плана) для решения ЗЛП состоит

- 1) умение находить начальный опорный план;
- 2) наличие признака оптимальности опорного плана;
- 3) умение переходить к нехудшему опорному плану.

Пусть ЗЛП представлена системой ограничений в каноническом виде:

$$\sum_{j=1}^n a_{ij}x_j = b_i, \quad b_i \geq 0 \quad (i = 1, \dots, m)$$

Говорят, что ограничение ЗЛП имеет предпочтительный вид, если при неотрицательной правой части ($b_i \geq 0$) левая часть ограничений содержит переменную, входящую с коэффициентом, равным единице, а в остальные ограничения равенства - с коэффициентом, равным нулю.

Пусть система ограничений имеет вид

$$\sum_{j=1}^n a_{ij}x_j \leq b_i, \quad b_i \geq 0 \quad (i = 1, \dots, m)$$

Сведем задачу к каноническому виду. Для этого прибавим к левым частям неравенств дополнительные переменные $x_{n+1} \geq 0$ ($i = 1, \dots, m$). Получим систему, эквивалентную исходной:

$$\sum_{j=1}^n a_{ij}x_j + x_{n+1} = b_i, \quad b_i \geq 0 \quad (i = 1, \dots, m)$$

$$\vec{x}_0 = (0; \underbrace{0, \dots, 0}_n; \underbrace{b_1, b_2, \dots, b_m}_m)$$

которая имеет предпочтительный вид

В целевую функцию дополнительные переменные вводятся с коэффициентами, равными нулю $c_{n+1} = 0$ ($i = 1, \dots, m$).

Пусть далее система ограничений имеет вид

$$\sum_{j=1}^n a_{ij} x_j \geq b_i, \quad b_i \geq 0 \quad (i = 1, \dots, m)$$

Сведём её к эквивалентной вычитанием дополнительных переменных $x_{n+1} \geq 0$ ($i = 1, \dots, m$) из левых частей неравенств системы. Получим систему

$$\sum_{j=1}^n a_{ij} x_j - x_{n+1} = b_i, \quad b_i \geq 0 \quad (i = 1, \dots, m)$$

Однако теперь система ограничений не имеет предпочтительного вида, так как дополнительные переменные x_{n+1} входят в левую часть (при $b_i \geq 0$) с коэффициентами, равными -1 .

$$\bar{x}_0 = (\underbrace{0, \dots, 0}_n; -b_1, -b_2, \dots, -b_m)$$

Поэтому, вообще говоря, базисный план не является допустимым. В этом случае вводится так называемый искусственный базис. К левым частям ограничений-равенств, не имеющих предпочтительного вида, добавляют искусственные переменные ω_i . В целевую функцию переменные ω_i , вводят с коэффициентом M в случае решения задачи на минимум и с коэффициентом $-M$ для задачи на максимум, где M - большое положительное число. Полученная задача называется M -задачей, соответствующей исходной. Она всегда имеет предпочтительный вид.

Пусть исходная ЗЛП имеет вид

$$\max(\min) Z = \sum_{j=1}^n c_j x_j, \quad (2.16)$$

$$\sum_{j=1}^n a_{ij} x_j = b_i, \quad b_i \geq 0 \quad (i = 1, \dots, m), \quad (2.17)$$

$$x_j \geq 0 \quad (j = 1, \dots, n), \quad (2.18)$$

причём ни одно из ограничений не имеет предпочтительной переменной. M -задача запишется так:

$$\max(\min) \bar{Z} = \sum_{j=1}^n c_j x_j - (+) \sum_{i=1}^m M \omega_i \quad (2.19)$$

$$\sum_{j=1}^n a_{ij} x_j + \omega_i = b_i, \quad (i = 1, \dots, m) \quad (2.20)$$

$$x_j \geq 0 \quad (j = \overline{1, n}), \quad \omega_i \geq 0, \quad (i = 1, \dots, m) \quad (2.21)$$

Задача (2.19) - (2.21) имеет предпочтительный план. Её начальный опорный план имеет вид

$$\bar{x}_0 = (\underbrace{0, 0, \dots, 0}_n; b_1, b_2, \dots, b_m)$$

Если некоторые из уравнений (2.17) имеют предпочтительный вид, то в них не следует вводить искусственные переменные.

Теорема. Если в оптимальном плане

$$\bar{x} = (x_1, x_2, \dots, x_n, \omega_1, \omega_2, \dots, \omega_m) \quad (2.22)$$

M -задачи (2.19) - (2.21) все искусственные переменные $\omega_i = 0$ ($i = 1, \dots, m$), то план $\bar{x} = (x_1, x_2, \dots, x_n)$ является оптимальным планом исходной задачи (2.16) - (2.18).

$z(\bar{x}^*) = f(\bar{y}^*)$. Если одна из двойственных задач неразрешима вследствие неограниченности целевой функции на множестве допустимых решений, то система ограничений другой задачи противоречива.

Теорема (об оценках). Двойственные оценки показывают приращение функции цели, вызванное малым изменением свободного члена соответствующего ограничения задачи математического программирования, точнее

$$\frac{\partial z(\bar{x}^*)}{\partial b_i} = y_i^* \quad (i = \overline{1, m})$$

1. Простейший поток, его свойства. Классификация потоков.

2. Марковские цепи, их свойства

Аппарат теории марковских процессов с дискретными состояниями и цепей Маркова широко используют в теории систем, в исследовании операций и других прикладных дисциплинах. Это обусловлено многими причинами, среди которых отметим следующие:

1) многие реальные технические системы имеют конечные множества возможных состояний, а их поведение в процессе функционирования адекватно моделируется марковскими процессами,

2) теория марковских процессов с дискретными состояниями и цепей Маркова разработана настолько глубоко, что позволяет решать широкий класс прикладных задач.

Марковские процессы Представление случайных процессов графом состояний

Рассмотрим физическую систему S , в которой протекает случайный процесс с дискретными состояниями: S_1, S_2, \dots, S_i , (1)

число которых конечно (или счетно). Состояния S_1, S_2, \dots могут быть качественными (т. е. описываться словами) или же каждое из них характеризуется случайной величиной (либо случайным вектором).

Прежде всего, рассмотрим множество состояний (1) с точки зрения его структуры - возможности системы S переходить из состояния s_j в данное состояние s_i - непосредственно или через другие состояния. Для этого удобно пользоваться наглядной схемой, так называемым графом состояний. Здесь и далее мы будем отчасти пользоваться терминологией теории графов. Имеется две основные разновидности графов: неориентированные и ориентированные.

Неориентированный граф - совокупность точек (вершин графа) с соединяющими некоторые из них отрезками (ребрами графа).

Ориентированный граф - это совокупность точек (вершин) с соединяющими некоторые из них ориентированными отрезками (стрелками).

При изложении теории случайных процессов с дискретными состояниями мы будем пользоваться только ориентированными графами. Вершины графа будут соответствовать состояниям системы. Вершину будем изображать прямоугольником, в который вписано обозначение состояния; стрелка, ведущая из вершины s_j в вершину s_i , будет обозначать возможность перехода системы S из состояния s_j в состояние s_i - непосредственно, минуя другие состояния. Стрелки графа могут изображаться не только прямолинейными, но и криволинейными отрезками (рис. 1). Сам граф системы S будем обозначать буквой G .

Переход по стрелке, ведущей из состояния s_i в него же, означает задержку системы в состоянии s_i . «Обратные стрелки» можно на графе не проставлять, так как все расчеты можно вести и без них.

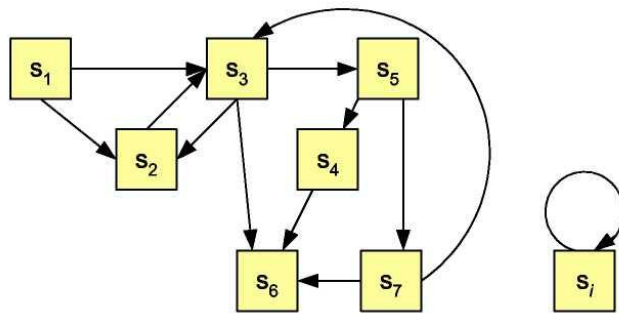


Рисунок 1 – Пример графа состояний

Проведем некоторую необходимую для дальнейшего классификацию состояний. Состояние s_i называется источником, если система S может выйти из этого состояния, но попасть в него обратно уже не может, т. е. на графе G состояний в состояние s_i не ведет ни одна стрелка. На рисунке 1 состояние s_1 является источником.

Состояние s_i называется конечным (или поглощающим), если система S может попасть в это состояние, но выйти из него уже не может. Для графа состояний это означает, что из состояния s_i не ведет ни одна стрелка (для графа, изображенного на рисунке 1, состояние s_6 поглощающее).

Если система S может непосредственно перейти из состояния s_i в состояние s_j то состояние s_j - называется соседним по отношению к состоянию s_i .

Состояние s_i называется транзитивным, если система S может войти в это состояние и выйти из него, т. е. на графе состояний есть хотя бы одна стрелка, ведущая в s_i и хотя бы одна стрелка, ведущая из s_i . На рисунке 1 все состояния, кроме s_1 и s_6 , являются транзитивными.

Для полноты картины можно рассматривать также и «изолированные» состояния. Состояние s_i называется изолированным, если из него нельзя попасть ни в одно из других состояний и в него нельзя попасть ни из какого другого состояния.

Наряду с отдельными состояниями системы S в ряде задач практически бывает нужно рассматривать подмножества ее состояний.

Обозначим W множество всех состояний системы S (конечное или бесконечное, но счетное) и рассмотрим его подмножество $V \subset W$. Подмножество V называется замкнутым (концевым), если система S , попав в одно (или находясь в одном) из состояний $s_i \in V$, не может выйти из этого подмножества состояний. Концевое подмножество состояний может включать в себя поглощающее состояние, а может и не включать.

Подмножество состояний $V \subset W$ называется связным или эргодическим, если из любого состояния, входящего в него, можно попасть в любое другое состояние, принадлежащее этому подмножеству. Эргодическим может быть и все множество W состояний системы S . В эргодическом множестве состояний нет ни источников, ни поглощающих состояний.

Подмножество состояний V называется транзитивным, если система S может войти в это подмножество и выйти из него, т. е. из любого состояния $s_i \in V$ можно (за то или другое число перескоков) выйти из этого подмножества.

Случайный процесс, протекающий в системе S , можно трактовать как процесс блуждания системы по множеству состояний W . Если подмножество $V \subset W$ является концевым, то, попав в него, система будет продолжать блуждание уже по этому подмножеству состояний V . Если все множество эргодично, то блуждание будет происходить по всем его состояниям.

На практике очень часто встречаются системы, состояния которых образуют цепь (рисунок 2), в которой каждое состояние s_i (кроме двух крайних s_0 и s_n) связано прямой и обратной связью с двумя соседними s_{i-1}, s_{i+1} , а каждое из двух крайних связано прямой и обратной связью только с одним соседним.

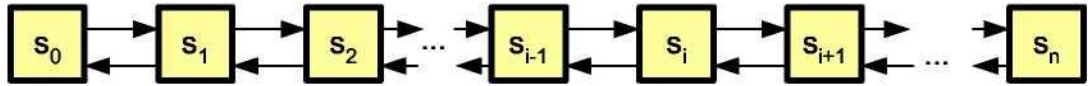


Рисунок 2 - Схема процесса гибели и размножения

Такая схема случайного процесса называется схемой гибели и размножения, а сам процесс — процессом гибели и размножения.

Если на графе состояний системы S стрелки, ведущие справа налево, отсутствуют, то говорят о процессе «чистого размножения», в противоположном случае — о процессе «чистой гибели».

Процесс гибели и размножения может в некоторых случаях иметь не конечное число состояний: $s_1, s_2, \dots, s_i, \dots, s_n$, а бесконечное (счетное): $s_1, s_2, \dots, s_i, \dots$.

При анализе случайных процессов, протекающих в системах с дискретными состояниями, важную роль играют вероятности состояний.

Обозначим $S(t)$ состояние системы S в момент t . Вероятностью i -го состояния в момент t называется вероятность события, состоящего в том, что в момент t система S будет в состоянии s_i . Обозначим ее $p_i(t)$:

$$p_i(t) = P\{S(t) = s_i\}, \quad (2)$$

где $S(t)$ - случайное состояние системы S в момент t . Очевидно, что для системы с дискретными состояниями $s_1, s_2, \dots, s_i, \dots$, в любой момент t сумма вероятностей состояний равна единице:

$$\sum_i p_i(t) = 1, \quad (3)$$

как сумма вероятностей полной группы несовместных событий.

В ряде задач практики нас интересует так называемый установившийся или стационарный режим работы системы, который в ней устанавливается, когда от начала процесса прошло достаточно большое время t . Например, процесс изменения напряжения в сети питания технического устройства, пройдя сразу после включения через ряд колебаний, по прошествии времени, устанавливается. Аналогично этому и в некоторых случайных процессах по прошествии достаточно большого времени t устанавливается стационарный режим, во время которого состояния системы хотя и меняются случайным образом, но их вероятности $p_i(t)$ ($i = 1, 2, \dots$) остаются постоянными. Обозначим эти постоянные вероятности p_i .

$$p_i = \lim p_i(t) \quad (4)$$

Вероятности p_i ($i = 1, 2, \dots$), если они существуют, называются финальными (предельными) вероятностями состояний. Финальную вероятность p_i можно истолковать как среднюю долю времени, которую в стационарном режиме проводит система S в состоянии s_i . В дальнейшем будет показано, при каких условиях финальные вероятности существуют и какими они могут быть для разных состояний и подмножеств состояний.

Введем очень важное для дальнейшего понятие марковского случайного процесса.

Случайный процесс, протекающий в системе S с дискретными состояниями $s_1, s_2, \dots, s_i, \dots$, называется марковским, если для любого момента времени t_0 вероятность каждого из состояний системы в будущем (при $t > t_0$) зависит только от ее состояния в настоящем (при $t = t_0$) и не зависит от того, когда и как она пришла в это состояние; т. е. не зависит от ее поведения в прошлом (при $t < t_0$).

Не надо понимать марковское свойство случайного процесса как полную независимость «будущего» от «прошлого»; в общем случае «будущее» зависит от «настоящего», т. е. вероятности $p_i(t)$ при $t > t_0$ зависят от того, в каком состоянии s_i находится система в настоящем (при $t=t_0$); само же это «настоящее» зависит от «прошлого», от того, как вела себя система S при $t < t_0$. Это можно сформулировать следующим образом: для марковского случайного процесса «будущее» зависит от «прошлого» только через «настоящее» (рисунок 3). При фиксированном «настоящем» условные вероятности всех состояний системы в «будущем» не зависят от предыстории процесса, т. е. от того, когда и как система S к моменту t_0 пришла в состояние

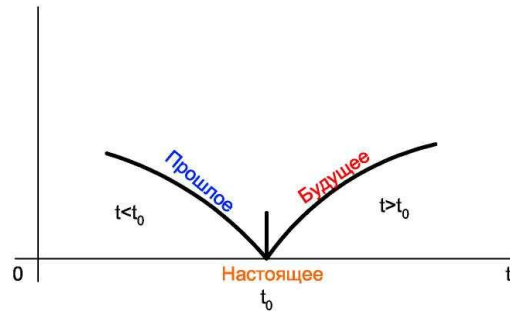


Рисунок 3 – Схема марковского свойства случайного процесса

«Настоящее» может быть задано не одним каким-то состоянием s_i , а целым подмножеством состояний $V \subset W$, где W - множество всех возможных состояний системы.

Подчеркнем также, что «настоящее» может быть задано не только одним состоянием системы S в момент t_0 ; в него при желании можно включить и те элементы из «прошлого», от которых, при заданном «настоящем», зависит будущее. Например, вероятности состояний в «будущем» могут зависеть не только от состояния s_i системы в настоящем, но и от того, из какого состояния s_i система перешла к моменту t_0 в состояние s_i ; в этом случае настоящее характеризуется не только состоянием s_i , в которое система перешла к моменту t_0 , но и состоянием s_j , из которого она перешла в s_i . Вводя в состав параметров, характеризующих настоящее состояние системы, те параметры из прошлого, от которых зависит будущее, можно, как говорится, «марковизировать» многие немарковские случайные процессы, но, как правило, это приводит к сильному усложнению математического аппарата.

3. Марковские процессы в инженерной практике

Марковские случайные процессы с дискретными состояниями и дискретным временем

Пусть имеется система S с дискретными состояниями $s_1, s_2, \dots, s_i, \dots, s_n$. Предположим, что случайные переходы («перескоки») системы из состояния в состояние могут происходить только в определенные моменты времени t_0, t_1, t_2, \dots . Эти моменты мы будем называть шагами процесса; $t_0=0$ - его началом. Сам процесс представляет собой случайное блуждание системы S по состояниям. После первого шага система может оказаться в одном (и только в одном) из своих возможных состояний: $s_1^{(1)}, s_2^{(1)}, \dots, s_i^{(1)}, \dots, s_n^{(1)}$; на втором шаге - $s_1^{(2)}, s_2^{(2)}, \dots, s_i^{(2)}, \dots, s_n^{(2)}$, на k -м шаге $s_1^{(k)}, s_2^{(k)}, \dots, s_i^{(k)}, \dots, s_n^{(k)}$ (число состояний в общем случае может быть бесконечным, но счетным. Здесь же для простоты ограничимся конечным числом n состояний).

Предположим, что граф состояний системы S имеет вид, представленный на рисунке 4. Процесс блуждания системы S по состояниям можно представить как последовательность или «цепь» событий, состоящих в том, что в начальный момент $t_0=0$ система находится в одном из состояний (например, в состоянии $s_1^{(0)}$), в момент первого

шага перешла из него скачком в состояние $s_5^{(1)}$, из которого на втором шаге перешла в $s_3^{(2)}$, на третьем шаге перешла в $s_2^{(3)}$ и т. д. «Траектория» системы, блуждающей по состояниям s_1, s_5, s_3, s_2 показана на рисунке 4 жирными линиями. На каких-то шагах система может задерживаться в том или другом из своих состояний, $s_i^{(k)} = s_i^{(k+1)}$ (это показано «возвратной стрелкой» на рисунке 4) или же вернуться в него после ряда шагов.

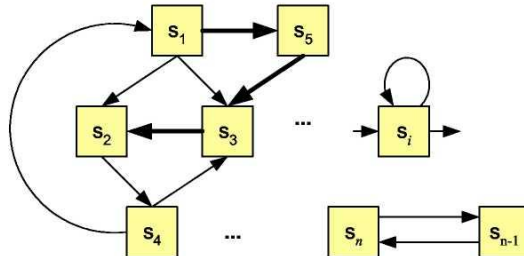


Рисунок 4 – Граф состояний системы S

«Траектория» блуждания системы по графу состояний, изображенная на рисунке 4 жирными линиями, представляет собой не что иное, как реализацию случайного процесса, полученную в результате одного опыта. При повторении опыта, естественно, реализации в общем случае не совпадают.

Рассмотрим общий случай. Пусть происходит случайный процесс в системе S с дискретными состояниями $s_1, s_2, \dots, s_i, \dots, s_n$, которые она может принимать в последовательности шагов с номерами $0, 1, 2, \dots, k, \dots$.

Случайный процесс представляет собой последовательность событий вида $\{S(k) = s_i\}$ ($i = 1, 2, \dots, n; k = 0, 1, 2, \dots$). Наиболее важной ее характеристикой являются вероятности состояний системы

$$P\{S(k) = s_i\} \quad (i = 1, 2, \dots, n; k = 0, 1, 2, \dots), \quad (5)$$

где $P\{S(k) = s_i\}$ - вероятность того, что на k -м шаге система S будет находиться в состоянии s_i .

Распределение вероятностей (5) представляет собой не что иное, как одномерный закон распределения случайного процесса $S(t)$, протекающего в системе S с «качественными» дискретными состояниями и дискретным временем $t_0, t_1, t_2, \dots, t_k$.

Процесс, протекающий в такой системе S , называется марковским процессом с дискретными состояниями и дискретным временем (или, короче, марковской цепью), если выполняется условие: для любого фиксированного момента времени (любого шага k_0) условные вероятности состояний системы в будущем (при $k > k_0$) зависят только от состояния системы в настоящем (при $k = k_0$) и не зависят от того, когда (на каком шаге, при $k < k_0$) и откуда система пришла в это состояние. Марковская цепь представляет собой разновидность марковского процесса, в котором будущее зависит от прошлого только через настоящее.

Цепь, в которой условные вероятности состояний в будущем зависят только от состояния на данном, последнем, шаге и не зависят от предыдущих, иногда называют простой цепью Маркова, в отличие от такой, где будущее зависит от состояний системы не только в настоящем на данном шаге, но и от ее состояний на нескольких предыдущих шагах; такую цепь называют сложной цепью Маркова. Сам А. А. Марков рассматривал сложные цепи, построенные на материале буквенных последовательностей, взятых из текста пушкинского «Евгения Онегина».

Если в качестве системы, в которой происходит случайный процесс, рассмотреть букву, входящую в текст, которой могут быть: а, б, в, щ, ь, ы, ь, э, ю, я, «пробел», то сразу ясно, что вероятность последующей буквы быть той или другой зависит от того, какова была предыдущая (например, последовательности букв «яы» или «эь» в русском языке исключены); не так очевидно, но все же ясно, что эта вероятность зависит не только от предыдущей буквы, но и от других, ей предшествовавших (например, последовательность букв «ттт» в русском языке если не исключена, то практически невозможна, тогда как последовательность «тт» встречается довольно часто). Мы в данном элементарном изложении будем рассматривать только простые цепи Маркова и вычислять для них вероятности состояний.

Из определения марковской цепи следует, что для нее вероятность перехода системы S в состояние s_i на $(k+1)$ -м шаге зависит только от того, в каком состоянии s_i находилась система на предыдущем k -м шаге и не зависит от того, как она вела себя до этого k -го шага.

Основной задачей исследования марковской цепи является нахождение безусловных вероятностей нахождения системы S на любом k -м шаге в состоянии s_i ; обозначим эту вероятность $p_i(k)$:

$$p_i(k) = P\{S(k) = s_i\} \quad (i = 1, 2, \dots, n; k = 0, 1, 2, \dots). \quad (6)$$

Для нахождения этих вероятностей необходимо знать условные вероятности перехода системы S на k -м шаге в состояние s_i , если известно, что на предыдущем $(k-1)$ -м шаге она была в состоянии s_j . Обозначим эту вероятность

$$p_{ij}(k) = P\{S(k) = s_i \mid S(k-1) = s_j\} \quad (i, j = 1, 2, \dots, n; k = 0, 1, 2, \dots). \quad (7)$$

Вероятности $p_{ij}(k)$ называются переходными вероятностями марковской цепи на k -м шаге. Вероятность $p_{ij}(k)$ есть вероятность того, что на k -м шаге система задержится (останется) в состоянии s_i .

Переходные вероятности $p_{ij}(k)$ можно записать в виде квадратной таблицы (матрицы) размерности $n \times n$:

$$\|p_{ij}(k)\| = \begin{pmatrix} p_{11}(k) & p_{12}(k) & \dots & p_{1j}(k) & \dots & p_{1n}(k) \\ p_{21}(k) & p_{22}(k) & \dots & p_{2j}(k) & \dots & p_{2n}(k) \\ \dots & \dots & \dots & \dots & \dots & \dots \\ p_{i1}(k) & p_{i2}(k) & \dots & p_{ij}(k) & \dots & p_{in}(k) \\ \dots & \dots & \dots & \dots & \dots & \dots \\ p_{n1}(k) & p_{n2}(k) & \dots & p_{nj}(k) & \dots & p_{nn}(k) \end{pmatrix} \quad (k = 0, 1, 2, \dots) \quad (8)$$

По главной диагонали матрицы (8) стоят вероятности задержки системы в данном состоянии s_j ($j = 1, \dots, n$) на k -м шаге.

$$p_{11}(k), p_{22}(k), \dots, p_{ii}(k), \dots, p_{nn}(k). \quad (9)$$

Так как на каждом шаге система S может находиться только в одном из взаимно исключающих состояний, то для любой k -й строки матрицы (8) сумма всех стоящих в ней вероятностей равна единице:

$$\sum_{j=1}^n p_{ij}(k) = 1. \quad (10)$$

Матрица, обладающая таким свойством, называется стохастической. Естественно, что все элементы стохастической матрицы отвечают условию $0 \leq p_{ij}(k) \leq 1$. В силу

условия (10) можно в матрице (8) не задавать вероятности задержки, а получать их как дополнения до единицы всех остальных членов строки:

$$p_{ii}(k) = 1 - \sum_{j=1}^n p_{ij}(k). \quad (11)$$

Чтобы найти безусловные вероятности $p_i(k)$, недостаточно знать матрицу переходных вероятностей (8); нужно еще знать начальное распределение вероятностей, т. е. вероятности состояний $p_i(0)$, соответствующие началу процесса - моменту $t_0 = 0$:

$$p_1(0), p_2(0), \dots, p_i(0), \dots, p_n(0), \quad (12)$$

в сумме образующие единицу:

$$\sum_{i=1}^n p_i(0) = 1 \quad (13)$$

Если известно, что в начальный момент система S находится во вполне определенном состоянии s_i , то вероятность $p_i(0)$ этого состояния в формуле (13) равна единице, а все остальные - нулю:

$$p_i(0), p_1(0) = p_2(0) = \dots = p_{i-1}(0) = p_{i+1}(0) = \dots = p_n(0) = 0. \quad (14)$$

Цепь Маркова называется однородной, если переходные вероятности $p_{ij}(k)$ не зависят от номера шага k : $p_{ij}(k) = p_{ij}$. Матрица переходных вероятностей для однородной цепи Маркова имеет вид:

$$\|p_{ij}\| = \begin{pmatrix} p_{11} & p_{12} \dots p_{1j} \dots p_{1n} \\ p_{21} & p_{22} \dots p_{2j} \dots p_{2n} \\ \dots & \dots & \dots & \dots \\ p_{i1} & p_{i2} \dots p_{ij} \dots p_{in} \\ \dots & \dots & \dots & \dots \\ p_{n1} & p_{n2} \dots p_{nj} \dots p_{nn} \end{pmatrix} \quad (14)$$

При выводе формул для вероятностей состояний, в целях простоты записи, будем рассматривать только однородные цепи Маркова (в случае, когда цепь неоднородна, можно все переходные вероятности в формулах просто положить зависящими от номера шага k).

При нахождении вероятностей состояний марковской цепи на k -м шаге $p_i(k)$ ($k = 1, 2, \dots$) удобно бывает пользоваться так называемым размеченным графом состояний системы S , где возле каждой стрелки, ведущей из состояния s_i в состояние s_j , проставлена переходная вероятность p_{ij} ; вероятности задержки на размеченном графе не проставляются, а просто получаются дополнением до единицы суммы вероятностей, стоящих у всех стрелок, ведущих из данного состояния s_i .

Теперь покажем, как найти для однородной цепи Маркова безусловную вероятность нахождения системы S на k -м шаге в состоянии s_j ($j = 1, 2, \dots, n$)

$$p_j(k) = \mathbf{P}\{S(k) = s_j\}, \quad (15)$$

если задана матрица переходных вероятностей $\|p_{ij}\|$ (или, что равнозначно, размеченный граф состояний) и начальное распределение вероятностей

$$p_i(0) \quad (i = 1, 2, \dots, n).$$

$$\sum_{j=1}^n p_{ij}(k) = 1 \quad (16)$$

Сделаем гипотезу, состоящую в том, что в начальный момент ($k=0$) система находилась в состоянии s_i . Вероятность этой гипотезы известна из (16) и равна

$p_i(0)=P\{S(0)=s_i\}$. В предположении, что эта гипотеза имеет место, условная вероятность того, что система S на первом шаге будет в состоянии s_j , равна переходной вероятности $p_{ij}(k) = \mathbf{P}\{S(1) = s_j \mid S(0) = s_i\}$.

По формуле полной вероятности получим:

$$p_j(1) = \sum_{i=1}^n P\{S(1) = s_j | S(0) = s_i\} \cdot P\{S(0) = s_i\} = \sum_{i=1}^n p_{ij} p_i(0), \quad (j = 1, \dots, n) \quad (17)$$

Таким образом, мы нашли распределение вероятностей системы S на первом шаге. Теперь у нас есть все необходимое для того, чтобы найти распределение вероятностей на втором шаге, которое для цепи Маркова зависит только от распределения вероятностей на первом шаге и матрицы переходных вероятностей.

Опять сделаем гипотезу, состоящую в том, что на первом шаге система находится в состоянии s_i вероятность этой гипотезы нам уже известна и равна $p_i(1) = P\{S(1) = s_i\}$. При этой гипотезе условная вероятность того, что на втором шаге система S будет в состоянии s_i , равна:

$$p_{ij}(k) = P\{S(2) = s_j | S(1) = s_i\}$$

По формуле полной вероятности находим

$$p_j(2) = \sum_{i=1}^n p_i(1) p_{ij}, \quad (j = 1, 2, \dots, n) \quad (18)$$

Таким образом, мы выразили распределение вероятностей (18) на втором шаге через распределение вероятностей на первом шаге и матрицу $\|p_{ij}\|$. Переходя таким же способом от $k = 2$ к $k = 3$ и т. д., получим рекуррентную формулу:

$$p_j(k) = \sum_{i=1}^n p_i(k-1) p_{ij}, \quad (k = 1, 2, \dots, n; j = 1, 2, \dots, n) \quad (19)$$

При некоторых условиях в цепи Маркова с возрастанием k (номера шага) устанавливается стационарный режим, в котором система S продолжает блуждать по состояниям, но вероятности этих состояний уже от номера шага не зависят. Такие вероятности называются предельными (или финальными) вероятностями цепи Маркова.

Например, если рассматривать ЭВМ в двух состояниях: s_1 - исправна, s_2 - не исправна, то имеет место следующая динамика изменения вероятностей (при начальных условиях):

$$p_1(0) = 1, p_2(0) = 0; p_1(1) = 0,7; p_1(2) = 0,61; p_1(3) = 0,583; p_1(4) = 0,5749.$$

Ниже мы покажем, что в этом случае $p_1 = \lim_{k \rightarrow \infty} p_1(k) = 0,4/(0,4+0,3) = 0,5714$. Таким образом, в рассматриваемой системе стационарный режим наступит практически через четыре шага.

Можно убедиться в том, что в этом примере финальные вероятности не зависят от начальных условий.

Сформулируем условия существования стационарного режима для системы S с конечным числом состояний n , в которой протекает марковский случайный процесс с дискретными состояниями и дискретным временем (цепь Маркова):

1. Множество всех состояний W системы S должно быть эргодическим.

2. Цепь Маркова должна быть однородной:

$$p_{ij}(k) = p_{ij} \quad (20)$$

3. Цепь Маркова должна быть «достаточно хорошо перемешиваемой» (не должна быть «циклической»).

Цепи Маркова, отвечающие этим условиям, будем называть эргодическими цепями Маркова.

1. 2 Лекция №7 (2 часа).

Тема: «Теоретические основы обработки экспериментальных данных»

1.2.1 Вопросы лекции:

1. Статистический материал и его первичная обработка. Эмпирические законы распределения. Полигон частот, гистограмма.
2. Числовые характеристики выборки. Точечные оценки выборочных характеристик. Интервальные оценки, их свойства.
3. Метод доверительных интервалов при заданных условиях. Метод моментов.
4. Статистические гипотезы, ошибки первого и второго рода.. Статистические критерии, их виды, мощность критерия.
5. Выравнивание статистических рядов.
6. Функция регрессии, коэффициент детерминации, корреляции, ковариация. Виды регрессий, статистическая значимость их параметров. Автокорреляция

1.2.2 Краткое содержание вопросов:

1. Статистический материал и его первичная обработка. Эмпирические законы распределения. Полигон частот, гистограмма.

Генеральная совокупность и выборка

Предметом математической статистики является изучение случайных величин (или случайных событий) по результатам наблюдений.

Для получения опытных данных необходимо провести обследование соответствующих объектов.

Совокупность всех мысленно возможных объектов данного вида, над которыми проводятся наблюдения с целью получения конкретных значений определённой случайной величины, называется **генеральной совокупностью**.

Генеральную совокупность будем называть **конечной** или **бесконечной** в зависимости от того, конечна или бесконечна совокупность составляющих её элементов.

Часть отобранных объектов из генеральной совокупности называется **выборочной совокупностью** или **выборкой**.

Число N объектов генеральной совокупности и число n объектов выборочной совокупности будем называть **объёмами генеральной и выборочной совокупности** соответственно.

Для того чтобы по выборке можно было достаточно уверенно судить о случайной величине, выборка должна быть **представительной (репрезентативной)**. Репрезентативность выборки означает, что объекты выборки достаточно хорошо представляют генеральную совокупность. Она обеспечивается случайностью отбора.

Существуют несколько способов отбора, обеспечивающих репрезентативность выборки. Рассмотрим некоторые из них.

После того как сделана выборка, все объекты этой совокупности обследуются по отношению к определённой случайной величине и получают наблюдаемые данные.

Для изучения закономерностей варьирования значений случайной величины опытные данные подвергают обработке.

Операция, заключающаяся в том, что результаты наблюдений над случайной величиной, т.е. наблюдаемые значения случайной величины, располагают в порядке неубывания, называется **ранжированием опытных данных**.

После операции ранжирования опытные данные объединяют в группы так, чтобы в каждой отдельной группе значения случайной величины будут одинаковы.

Значение случайной величины, соответствующее отдельной группе сгруппированного ряда наблюдаемых данных, называется вариантом (x_i) (**вариантой**), а изменение этого значения – **варьированием**.

Численность отдельной группы сгруппированного ряда наблюдаемых данных называется **частотой** или **весом** (m_i) соответствующей **варианты**.

Отношение частоты данного варианта к общей сумме частот всех вариантов

называется **частотой** или долей этой **варианты** (P_i):
$$P_i = \frac{m_i}{\sum_{i=1}^v m_i},$$

где v – число вариантов. Полагая $n = \sum_{i=1}^v m_i$, где n – объём выборки, имеем: $P_i = \frac{m_i}{n}$.

Заметим, что частота P_i – статистическая вероятность появления варианта x_i .

Дискретным вариационным рядом распределения называется ранжированная совокупность вариантов x_i , с соответствующими им частотами m_i или частотами P_i .

Если изучаемая случайная величина является непрерывной, то ранжирование и группировка наблюдаемых значений зачастую не позволяют выделить характерные черты варьирования её значений. Это объясняется тем, что отдельные значения случайной величины могут как угодно мало отличаться друг от друга и поэтому в совокупности наблюдаемых данных одинаковые значения случайной величины могут встречаться редко, а частоты вариантов мало отличаются друг от друга.

Интервальным вариационным рядом называется упорядоченная совокупность интервалов варьирования значений случайной величины с соответствующими частотами или частотами попаданий в каждый из них значений величины.

Рассмотрим алгоритм построения интервального ряда.

1. Для построения интервального ряда необходимо определить величину частичных интервалов, на которые разбивается весь интервал варьирования наблюдаемых значений случайной величины. Считая, что все частичные интервалы имеют одну и ту же длину, для каждого интервала следует установить его верхнюю и нижнюю границы, а затем в соответствии с полученной упорядоченной совокупностью частичных интервалов сгруппировать результаты наблюдений. Длину частичного интервала h следует выбрать так, чтобы построенный ряд не был громоздким и в то же время позволил выявить характерные черты изменения значений случайной величины.

2. Найдём размах варьирования ряда R : $R = x_{\text{наиб}} - x_{\text{наим}}$

Выберем число интервалов v (обычно от 7 до 11).

3. Для более точного определения величины частичного интервала можно

воспользоваться **формулой Стерджеса**:
$$h = \frac{R}{1 + 3,322 \lg n}.$$

Если h – дробное, то за длину частичного интервала следует брать ближайшее целое число, либо ближайшую простую дробь.

4. За начало первого интервала следует брать величину: $x_{нач} = x_{наим} - 0,5h$.

5. Конец последнего интервала ($x_{кон}$) должен удовлетворить условию: $x_{кон} - h \leq x_{наиб} < x_{кон}$.

6. Промежуточные интервалы получают, прибавляя к концу предыдущего интервала длину частичного интервала h .

7. Определим, сколько значений признака попало в каждый конкретный интервал. При этом в интервал включают значения случайной величины, большие или равные нижней границе и меньшие верхней границы. Иногда интервальный вариационный ряд для простоты исследования условно заменяют дискретным. В этом случае серединное значение i -го интервала принимают за вариант x_i , а соответствующую интервальную частоту m_i – за частоту этой варианты.

2. Числовые характеристики выборки.

Закон распределения (или просто распределение) случайной величины можно задать различными способами. Например, дискретную случайную величину можно задать с помощью или ряда распределения, или интегральной функции, а непрерывную случайную величину – с помощью или интегральной, или дифференциальной функции. Рассмотрим выборочные аналоги этих двух функций.

В теории вероятностей для характеристики распределения случайной величины X служит интегральная функция распределения $F(x) = P(X < x)$. В дальнейшем, если величина X распределена по некоторому закону $F(x)$, будем говорить, что и генеральная совокупность распределена по закону $F(x)$. Введём выборочный аналог функции $F(x)$.

Пусть имеется выборочная совокупность значений некоторой случайной величины X объёма n и каждому варианту из этой совокупности поставлена в соответствие его частота. Пусть, далее, x – некоторое действительное число, а m_x – число выборочных значений случайной величины X , меньших x . Тогда число m_x/n является частотой наблюдаемых в выборке значений величины X , меньших x , т.е. частотой появления события $X < x$. При изменении x в общем случае будет изменяться и величина m_x/n . Это означает, что относительная частота m_x/n является функцией аргумента x . А т.к. эта функция находится по выборочным данным, полученным в результате опытов, то её называют **выборочной** или **эмпирической**.

Выборочной функцией распределения (или **функцией распределения выборки**) называется функция $F(x)^*$, задающая для каждого значения x относительную частоту события $X < x$.

Итак, по определению, $F(x)^* = m_x/n$, где n – объём выборки, m_x – число выборочных значений случайной величины X , меньших x . В отличие от выборочной функции $F(x)^*$ интегральную функцию $F(x)$ генеральной совокупности называют **теоретической функцией распределения**. Главное различие функций $F(x)$ и $F(x)^*$ состоит в том, что теоретическая функция распределения $F(x)$ определяет вероятность события $X < x$, а выборочная функция – относительную частоту этого события.

Свойство статистической устойчивости частоты, обоснованное теоремой Бернулли, оправдывает целесообразность использования функции $F(x)^*$ при больших n в качестве приближённого значения неизвестной функции $F(x)$.

В заключение отметим, что функция $F(x)$ и её выборочный аналог $F(x)^*$ обладают одинаковыми свойствами. Действительно, из определения функции $F(x)^*$ имеем следующие свойства:

1. $0 \leq F^*(x) \leq 1$ 2. $F^*(x)$ – неубывающая функция. 3. $F^*(-\infty) = 0, F(\infty) = 1$.

Таковыми же свойствами обладает и функция $F(x)$.

Наблюдаемые данные, представленные в виде вариационного ряда, можно изобразить графически, используя не только функцию $F^*(x)$. К наиболее распространённым видам графического изображения вариационных рядов относятся **полигон** и **гистограмма**. Графическое изображение рядов с помощью полигона или гистограммы позволяет получить наглядное представление о закономерности варьирования наблюдаемых значений случайной величины.

Полигон обычно используют для изображения дискретного вариационного ряда. Для его построения в прямоугольной системе координат наносят точки с координатами $(x_i; m_i)$ или $(x_i; p^*_i)$, где x_i – значение i -го варианта, а m_i (p^*_i) – соответствующие частоты (частоты). Затем отмеченные точки соединяют отрезками прямой линии. Полученная ломаная называется **полигоном**.

Если полигон частот построен по дискретному вариационному ряду дискретной случайной величины, то его называют **многоугольником распределения частот**, который является выборочным аналогом многоугольника распределения вероятностей. Заметим, что сумма ординат многоугольника распределения частот, как и у многоугольника распределения вероятностей, равна 1, т.к. $\sum p^*_i = 1$.

Гистограмма служит только для изображения интервальных вариационных рядов. Для её построения в прямоугольной системе координат на оси Ox откладывают отрезки, изображающие частичные интервалы варьирования, и на этих отрезках, как на основаниях, строят прямоугольники с высотами, равными частотам или частостям соответствующих интервалов. В результате такой операции получают ступенчатую фигуру, состоящую из прямоугольников, которую называют **гистограммой**.

Для графического изображения интервального вариационного ряда можно использовать полигон, если этот ряд преобразовать в дискретный. В этом случае интервалы заменяют их серединными значениями и ставят им в соответствие интервальные частоты (частоты). Для полученного дискретного ряда строят полигон.

Построив вариационный ряд и изобразив его графически, можно получить первоначальное представление о закономерностях, имеющих место в ряду наблюдений. Однако на практике зачастую этого недостаточно. Такая ситуация возникает, когда следует уточнить те или иные сведения о ряде распределения или когда имеется необходимость сравнить два ряда и более. При этом следует сравнивать однотипные вариационные ряды, т.е. такие ряды, которые получены при обработке сравнимых статистических данных.

Сравниваемые распределения могут существенно отличаться друг от друга. Они могут иметь различные средние значения случайной величины, вокруг которых группируются в основном остальные значения, или различаться рассеиванием данных наблюдений вокруг указанных значений и т.д. Поэтому для дальнейшего изучения изменения значений случайной величины используют числовые характеристики вариационных рядов. Поскольку эти характеристики вычисляются по статистическим данным (данным,

полученным в результате наблюдений), их обычно называют **статистическими характеристиками** или **оценками**.

Пусть собранный и обработанный статистический материал представлен в виде вариационного ряда. Теперь результаты наблюдений над случайной величиной следует подвергнуть анализу и выявить характерные особенности поведения случайной величины. Для этого удобнее всего выделить некоторые постоянные, которые представляли бы вариационный ряд в целом и отражали присущие изучаемой совокупности закономерности.

Некоторые из этих постоянных отличаются тем, что вокруг них концентрируются остальные результаты наблюдений. Такие величины называются **средними величинами**. К ним относятся среднее арифметическое (среднее выборочное), среднее геометрическое, среднее гармоническое и т.д. Однако эти характеристики не отражают «величину изменчивости» наблюдаемых данных, например величину разброса значений признака вокруг среднего арифметического. Другими словами, упомянутые средние величины не отражают вариацию.

Для характеристики изменчивости случайной величины, т.е. вариации, служат показатели вариации. К ним относятся размах варьирования R , среднее квадратическое отклонение, дисперсия и т.д.

2. Точечные оценки выборочных характеристик. Интервальные оценки, их свойства. Метод доверительных интервалов при заданных условиях. Метод моментов.

Выборочная характеристика, используемая в качестве приближённого значения неизвестной генеральной характеристики, называется её **точечной статистической оценкой**.

Среднее арифметическое \bar{O} – это точечная статистическая оценка математического ожидания $M(X)$; $D^*(X)$ – оценка дисперсии $D(X)$.

«Точечная» означает, что оценка представляет собой число или точку на числовой оси. «Статистическая» означает, что оценка рассчитывается по результатам наблюдений, т.е. по собранной исследователем статистике. Далее слово «статистическая» будет опускаться.

Обозначим через Θ («тэта») некоторую генеральную характеристику (ею может быть и MX , и любая другая числовая характеристика случайной величины X). Её числовое значение неизвестно, однако предложен некоторый алгоритм или формула вычисления точечной оценки $\Theta_{(n)}$ этой характеристики по результатам X_1, X_2, \dots, X_n наблюдений величины X . Обозначая буквой f этот алгоритм, запишем $\Theta^*_{(n)} = f(X_1, X_2, \dots, X_n)$. (3)

Подставив в (3) вместо X_1, X_2, \dots, X_n конкретные результаты наблюдений (конкретные числа), получим число, которое и принимают за приближённое значение неизвестной генеральной характеристики Θ . Найти погрешность этого приближения нельзя, поскольку числовое значение характеристики Θ неизвестно. Чтобы ответить на вопрос, хорошо или нет найденное приближение, рассмотрим оценку $\Theta^*_{(n)}$ с других позиций.

Пусть в формуле (3) X_1, X_2, \dots, X_n – не конкретные числа, а лишь обозначения тех результатов наблюдений, которые мы хотели бы получить. Но результат каждого отдельного наблюдения случайной величины случаен, т.е. X_1, X_2, \dots, X_n – это случайные величины, поэтому и оценка $\Theta^*_{(n)}$ также величина случайная; следовательно, можно говорить о её математическом ожидании ($M(\Theta^*_{(n)})$), дисперсии ($D(\Theta^*_{(n)})$) и законе распределения. Интерпретация оценки $\Theta^*_{(n)}$ как случайной величины позволяет

сформулировать свойства, которыми должна была обладать оценка, чтобы её можно было считать хорошим приближением к неизвестной генеральной характеристике. Это свойства состоятельности, несмещённости и эффективности.

Оценка $\Theta^*_{(n)}$ генеральной характеристики Θ называется **состоятельной**, если для любого $\varepsilon > 0$ выполняется равенство $\lim_{n \rightarrow \infty} P(|\Theta^*_{(n)} - \Theta| < \varepsilon) = 1$. (4)

Поясним смысл равенства (4). Пусть ε – очень малое положительное число. Тогда равенство (4) означает, что чем больше число наблюдений n , тем больше уверенность (вероятность) в незначительном по абсолютной величине отклонении оценки $\Theta^*_{(n)}$ от неизвестной характеристики Θ или короче: чем больше объём исходной информации, тем «ближе мы к истине». Если это так, то $\Theta^*_{(n)}$ – состоятельная оценка.

«Хорошая» оценка обязательно должна обладать свойством состоятельности. В противном случае оценка не имеет практического смысла: увеличение объёма исходной информации не будет «приближать нас к истине». Поэтому свойство состоятельности следует проверять в первую очередь.

Оценка $\Theta^*_{(n)}$ генеральной характеристики Θ называется **несмещённой**, если для любого фиксированного числа наблюдений n выполняется равенство $M(\Theta^*_{(n)}) = \Theta$, (5) т.е. математическое ожидание оценки равно неизвестной характеристике.

Несмещённая оценка $\Theta^*_{(n)}$ характеристики Θ называется **несмещённой эффективной**, если она среди всех прочих несмещённых оценок той же самой характеристики обладает наименьшей дисперсией.

Интервальные оценки параметров статистического распределения. Доверительные вероятности

Вычисляя на основании результатов наблюдений точечную характеристику Θ^* неизвестной числовой характеристики Θ , мы понимаем, что величина Θ^* является лишь приближённым значением характеристики Θ . Если для большого числа наблюдений точность приближения бывает достаточной для практических выводов (в силу несмещённости, состоятельности и эффективности «хороших» оценок), то для выборок небольшого объёма вопрос о точности оценок очень важен. В математической статистике он решается следующим образом. По сделанной выборке находится точечная оценка Θ^* неизвестной характеристики Θ , затем задаются вероятностью γ и по определённым правилам находят такое число $\varepsilon > 0$, чтобы выполнялось соотношение

$$P(\Theta^* - \varepsilon < \Theta < \Theta^* + \varepsilon) = \gamma. \quad (8)$$

Соотношению (8) тождественно соотношению

$$P(|\Theta^* - \Theta| < \varepsilon) = \gamma, \quad (9)$$

из которого видно, что абсолютная погрешность оценки Θ не превосходит числа ε . Это верно с вероятностью, равной γ . Число ε называется **точностью оценки Θ^*** (чем меньше ε , тем выше точность оценки), числа Θ_1 и Θ_2 называются **доверительными границами**, интервал (Θ_1, Θ_2) – **доверительным интервалом** или **интервальной оценкой** характеристики Θ , вероятность γ называется **доверительной вероятностью** или **надёжностью** интервальной оценки.

В соотношении (8) случайными величинами являются доверительные границы Θ_1 и Θ_2 : во-первых, эти границы могут изменяться при переходе от одной выборки к другой хотя бы потому, что при этом изменяется значение оценки Θ^* ; во-вторых, при фиксированной выборке границы Θ_1 и Θ_2 изменяются при изменении вероятности γ ,

поскольку ε выбирается в зависимости от γ . Генеральная же характеристики Θ – постоянная величина. Поэтому соотношение (8) следует читать так: «вероятность того, что интервал (Θ_1, Θ_2) накроет характеристику Θ , равна γ »; именно «интервал накроет характеристику», а не «характеристика попадёт в интервал».

Надёжность γ принято выбирать равной 0,95; 0,99; 0,999. Тогда событие, состоящее в том, что интервал (Θ_1, Θ_2) накроет характеристику Θ , будет практически достоверным. Также практически достоверным является событие, состоящее в том, что погрешность оценки Θ^* меньше ε , или, иначе, точность оценки Θ^* больше ε .

В соотношении (8) границы Θ_1 и Θ_2 симметричны относительно точечной оценки Θ^* . Обратим внимание на то, что не всегда удаётся построить границы с таким свойством.

Поскольку довольно часто встречаются нормально распределённые случайные величины, построим интервальные оценки для параметров нормального распределения – математического ожидания a и среднего квадратического отклонения σ .

4. Статистические гипотезы, ошибки первого и второго рода. Статистические критерии, их виды, мощность критерия.

Под **статистической гипотезой** понимают всякое высказывание о генеральной совокупности (случайной величине), проверяемое по выборке (по результатам наблюдений). Примером статистических гипотез являются следующие высказывания: генеральная совокупность, о которой мы располагаем лишь выборочными сведениями, имеет нормальный закон распределения или генеральная средняя (математическое ожидание случайной величины) равна 5. Не располагая сведениями о всей генеральной совокупности, высказанную гипотезу сопоставляют, по определённым правилам, с выборочными сведениями, и делают вывод о том, можно принять гипотезу или нет. Процедура сопоставления высказанной гипотезы с выборочными данными называется **проверкой гипотезы**.

Рассмотрим этапы проверки гипотезы и используемые при этом понятия.

Этап 1. Располагая выборочными данными X_1, X_2, \dots, X_n и руководствуясь конкретными условиями рассматриваемой задачи, формулируют гипотезу H_0 , которую называют **основной** или **нулевой**, и гипотезу H_1 , **конкурирующую** с гипотезой H_0 .

Термин «конкурирующая» означает, что являются противоположными следующие два события:

- по выборке будет принято решение о справедливости для генеральной совокупности гипотезы H_0 ;
- по выборке будет принято решение о справедливости для генеральной совокупности гипотезы H_1 .

Гипотезу H_1 называют также **альтернативной**.

Например, если нулевая гипотеза такова: математическое ожидание равно 5, – то альтернативная гипотеза может быть следующей: математическое ожидание меньше 5, что записывается следующим образом: $H_0 : M(X) = 5$; $H_1 : M(X) < 5$.

Этап 2. Задаются вероятностью α («альфа»), которую называют **уровнем значимости**. Поясним её смысл:

Решение о том, можно ли считать высказывание H_0 справедливым для генеральной совокупности, принимается по выборочным данным, т.е. по ограниченному ряду наблюдений; следовательно, это решение может быть ошибочным. При этом может иметь место ошибка двух родов:

- отвергают гипотезу H_0 , или, иначе, принимают альтернативную гипотезу H_1 , тогда как на самом деле гипотеза H_0 верна – это **ошибка первого рода**;

- принимают гипотезу H_0 , тогда как на самом деле высказывание H_0 неверно, т.е. верной является гипотеза H_1 – это **ошибка второго рода**.

Так вот, уровень значимости α – это вероятность ошибки первого рода, т.е. $\alpha = D_{I_0}(H_1)$, (13)

где $D_{I_0}(H_1)$ – вероятность того, что будет принята гипотеза H_1 , если на самом деле в генеральной совокупности верна гипотеза H_0 . Вероятность α задаётся заранее, разумеется, малым числом, поскольку это вероятность ошибочного заключения, при этом обычно используют некоторые стандартные значения: 0,05; 0,01; 0,005; 0,001. Например, $\alpha = 0,05$ означает следующее: если гипотезу H_0 проверять по каждой из 100 выборок одинакового объёма, то в среднем в 5 случаях из 100 мы совершим ошибку первого рода.

Вероятность ошибки второго рода обозначают β , т.е. $\beta = P_{H_1}(H_0)$, (14)

где $P_{H_1}(H_0)$ – вероятность того, что будет принята гипотеза H_0 , если на самом деле верна гипотеза H_1 . Зная α , можно найти вероятность β .

Обратим внимание на то, что в результате проверки гипотезы относительно гипотезы H_0 может быть принято и правильное решение. Существует правильное решение двух следующих видов:

- принимают гипотезу H_0 , тогда как и в действительности, в генеральной совокупности, она имеет место; вероятность этого решения $P_{H_0}(H_0) = 1 - \alpha$;

- не принимают гипотезу H_0 (т.е. принимают гипотезу H_1), тогда как на самом деле гипотеза H_0 неверна (т.е. верна гипотеза H_1); вероятность этого решения $P_{H_1}(H_1) = 1 - \beta$.

Этап 3. Находят величину ϕ такую, что:

- её значения зависят от выборочных данных X_1, X_2, \dots, X_n , т.е. для которой справедливо равенство $\phi = \phi(X_1, X_2, \dots, X_n)$;

- её значения позволяют судить о «расхождении выборки с гипотезой H_0 »;

- и она, будучи величиной случайной в силу случайности выборки X_1, X_2, \dots, X_n , подчиняется при выполнении гипотезы H_0 некоторому известному, затабулированному закону распределения.

Статистический критерий. Мощность критерия

Величину ϕ называют **критерием**.

Отметим, что в основе метода построения критерия лежит понятие функции правдоподобия.

Этап 4. Далее рассуждают так. Т.к. значения критерия позволяют судить о «расхождении выборки с гипотезой H_0 », то из области допустимых значений критерия ϕ следует выделить подобласть ω таких значений, которые свидетельствовали бы о существенном расхождении выборки с гипотезой H_0 , и, следовательно, о невозможности принять гипотезу H_0 . Подобласть ω называют **критической областью**. Допустим, что критическая область выделена. Тогда руководствуются следующим правилом: если вычисленное по выборке значение критерия ϕ попадает в критическую область, то гипотеза H_0 отвергается

и принимается гипотеза H_1 . При этом следует понимать, что такое решение может оказаться ошибочным: на самом деле гипотеза H_0 может быть справедливой. Т.обр., ориентируясь на критическую область, можно совершить ошибку первого рода, вероятность которой задана заранее и равна α . Отсюда вытекает следующее требование к критической области ω :

вероятность того, что критерий ϕ примет значение из критической области ω , должна быть равна заданному числу α , т.е.

$$P(\phi \in \omega) = \alpha. \quad (15)$$

Однако критическая область равенством (15) определяется неоднозначно. Действительно, представив себе график функции плотности $f_\phi(x)$ критерия ϕ , нетрудно понять, что на оси абсцисс существует бесчисленное множество областей-интервалов таких, что площади построенных на них криволинейных трапеций равны α , т.е. областей, удовлетворяющих требованию (15). Поэтому кроме требования (15) выдвигается следующее требование: критическая область ω должна быть расположена так, чтобы при заданной вероятности α ошибки первого рода вероятность β ошибки второго рода была минимальной.

Возможны три вида расположения критической области (в зависимости от вида нулевой и альтернативной гипотез, вида и расположения критерия ϕ):

правосторонняя критическая область, состоящая из интервала $(x_{\text{пр}, \alpha}^{\text{кр}}, +\infty)$, где точка $x_{\text{пр}, \alpha}^{\text{кр}}$ определяется из условия

$$P(\phi > x_{\text{пр}, \alpha}^{\text{кр}}) = \alpha \quad (16)$$

и называется **правосторонней критической точкой**, отвечающей уровню значимости α ;

левосторонняя критическая область, состоящая из интервала $(-\infty, x_{\text{лев}, \alpha}^{\text{кр}})$, где точка $x_{\text{лев}, \alpha}^{\text{кр}}$ определяется из условия

$$P(\phi < x_{\text{лев}, \alpha}^{\text{кр}}) = \alpha \quad (17)$$

и называется **левосторонней критической точкой**, отвечающей уровню значимости α ;

двусторонняя критическая область, состоящая из следующих двух интервалов:

$$(-\infty, x_{\text{лев}, \alpha/2}^{\text{кр}})$$

По значению критерия ϕ судят о «расхождении выборочных данных с гипотезой H_0 ». Естественно, что гипотеза H_0 должна быть отвергнута, если расхождения велики; именно этим объясняется включение в критическую область больших значений критерия ϕ (больше, чем критическая точка).

Включение же в ряде случаев в критическую область малых значений критерия ϕ (меньше, чем критическая точка) на первый взгляд противоречит смыслу этой величины. Однако не следует забывать, что ϕ – случайная величина (она зависит от результатов наблюдений X_1, X_2, \dots, X_n , которые случайны), поэтому маловероятно появление не только слишком больших, но и слишком малых её значений и их следует включить в критическую область.

Этап 5. В формулу критерия $\phi = \phi(X_1, X_2, \dots, X_n)$ вместо X_1, X_2, \dots, X_n подставляют конкретные числа, полученные в результате наблюдений, и подсчитывают числовое значение $\phi_{\text{чис}}$ критерия.

Если $\phi_{\text{чис}}$ попадает в критическую область ω , то гипотеза H_0 отвергается и принимается гипотеза H_1 . Поступая таким образом, следует понимать, что можно допустить ошибку с вероятностью α .

Если $\phi_{\text{чис}}$ не попадает в критическую область, гипотеза H_0 не отвергается. Но это вовсе не означает, что H_0 является единственно подходящей гипотезой: просто расхождение

между выборочными данными и гипотезой H_0 невелико, или, иначе, H_0 не противоречит результатам наблюдений; однако таким же свойством наряду с H_0 могут обладать и другие гипотезы.

5. Выравнивание статистических рядов. Критерий согласия Пирсона

Выше рассматривались гипотезы, относящиеся к отдельным параметрам распределения случайных величин, причём модели законов распределения этих величин представлялись известными. Однако во многих практических задачах модель закона распределения заранее не известна и возникает задача выбора модели, согласующейся с результатами наблюдений над случайной величиной.

Пусть высказано предположение, что неизвестная функция распределения $F_X(x)$ исследуемой случайной величины X имеет вполне определённую модель $F_{\text{теор}}(x)$, т.е. высказана гипотеза

$$H_0 : F_X(x) = F_{\text{теор}}(x). \quad (18)$$

В качестве теоретической модели $F_{\text{теор}}(x)$ может быть рассмотрена нормальная, биномиальная или какая-либо другая модель. Это определяется сущностью изучаемого явления, а также результатом предварительной обработки наблюдений над случайной величиной (формой графика вариационного ряда, соотношениями между выборочными характеристиками и т.д.).

Критерии, с помощью которых проверяется гипотеза (19), называются **критериями согласия**. Рассмотрим лишь один из них, использующий χ^2 -распределение и получивший название **критерия согласия Пирсона**.

Критерий предполагает, что результаты наблюдений сгруппированы в вариационный ряд. Для определённости положим, что это дискретный вариационный ряд с числом групп, равным v (см. строки 1 и 2 табл. 7).

Таблица 7.

x_i	x_1	...	x_{v-1}	x_v
m_i	m_1	...	m_{v-1}	m_v
$p_i^{\text{теор}} = P(X = x_i)$	$p_1^{\text{теор}} = P(X = x_1)$...	$p_{v-1}^{\text{теор}} = P(X = x_{v-1})$	$p_v^{\text{теор}} = 1 - p_1^{\text{теор}} - \dots - p_{v-1}^{\text{теор}}$
$m_i^{\text{теор}} = np_i^{\text{теор}}$	$m_1^{\text{теор}} = np_1^{\text{теор}}$...	$m_{v-1}^{\text{теор}} = np_{v-1}^{\text{теор}}$	$m_v^{\text{теор}} = np_v^{\text{теор}}$

Однако, прежде чем рассматривать сам критерий Пирсона, вспомним параметрическое оценивание закона распределения. Последовательность оценивания такая: формулируют гипотезу о модели закона распределения случайной величины; по результатам наблюдений находят оценки неизвестных параметров этой модели (допустим, что число неизвестных параметров равно l); вместо неизвестных параметров подставляют в модель найденные оценки. В результате предполагаемая модель закона оказывается полностью определённой и, используя её, рассчитывают вероятности $p_i^{\text{теор}} = P(X = x_i)$ того, что случайная величина X примет зафиксированные в наблюдениях значения x_i , $i=1, 2, \dots, v-1$; эти вероятности называют **теоретическими**. Обратим внимание на следующее обстоятельство: т.к. сумма вероятностей ряда распределения должна быть равна единице, т.е.

$$\sum_i p_i^{\text{теор}} = 1, \quad (19)$$

то полагаем вероятность $p_v^{\text{теор}} = 1 - p_1^{\text{теор}} - p_2^{\text{теор}} - \dots - p_{v-1}^{\text{теор}}$. Теоретические вероятности записаны в строке 3 табл. 7. Теперь найдём теоретические частоты $m_i^{\text{теор}} = np_i^{\text{теор}}$; они записаны в строке 4 табл. 7.

Обратим внимание на следующее: критерий согласия Пирсона можно использовать только в том случае, когда

$$m_i^{\text{теор}} \geq 5, i=1, 2, \dots, v. \quad (20)$$

Поэтому ту группу вариационного ряда, для которой это условие не выполняется, объединяют с соседней и соответственно уменьшают число групп; так поступают до тех пор, пока для каждой новой группы $m_i^{\text{теор}}$ будет не меньше 5. Новое число групп, как и прежде, обозначим символом v .

Оказывается, что если предполагаемая модель закона распределения действительно имеет место, т.е. верна гипотеза (18), и если к тому же выполняются условия (19) и (20), то величина

$$\varphi = \sum_{i=1}^v \frac{(m_i - m_i^{\text{теор}})^2}{m_i^{\text{теор}}} \quad (21)$$

будет иметь χ^2 -распределение с числом степеней свободы $k = v - l - 1$, т.е.

$$\varphi = \sum_{i=1}^v \frac{(m_i - m_i^{\partial \hat{a} \hat{b} \hat{c}})^2}{m_i^{\partial \hat{a} \hat{b} \hat{c}}} = \chi^2(k = v - l - 1),$$

где v – число (новое) групп вариационного ряда; l – число неизвестных параметров предполагаемой модели, оцениваемых по результатам наблюдений (если все параметры предполагаемого закона известны точно, то $l = 0$). Величину (21) и называют **критерием согласия χ^2** или **критерием согласия Пирсона**.

Далее поступаем так же, как обычно при проверке гипотез. Задаёмся уровнем значимости α . Зная распределение критерия φ , находим критическую область, как правило, это область правосторонняя, т.е. она имеет вид $(x_{\text{кр}, \alpha}^{\text{кр}}, +\infty)$; найдём числовое значение $\varphi_{\text{чис}}$ критерия (21). Если $\varphi_{\text{чис}}$ попадает в интервал $(x_{\text{кр}, \alpha}^{\text{кр}}, +\infty)$, то делаем вывод о неправомерности гипотезы H_0 (18); при этом не следует забывать, что этот вывод может оказаться ошибочным (на самом деле в генеральной совокупности гипотеза H_0 (18) имеет место) и вероятность того, что вывод ошибочен, равна α .

Если $\varphi_{\text{чис}}$ не попадает в интервал $(x_{\text{кр}, \alpha}^{\text{кр}}, +\infty)$, то гипотеза H_0 (18) не отвергается.

В заключение приведём схему определения точки $x_{\text{кр}, \alpha}^{\text{кр}}$:

$$\left. \begin{array}{l} \alpha \rightarrow \gamma = 1 - \alpha \\ l, v \rightarrow k = v - l - 1 \end{array} \right\} \longrightarrow \chi_{\gamma}^2 \rightarrow x_{\text{кр}, \alpha}^{\text{кр}} = \chi_{\gamma}^2. \quad (22)$$

Выборочный коэффициент корреляции

Выясним, можно ли измерить степень корреляционной и стохастической зависимости величины Y от X . Ответ проиллюстрируем примером 5. Все полученные в примере результаты объединены в табл. 8.

Таблица 8.

x_i	$x_1 = 2$	$x_2 = 5$	$x_3 = 8$
$P(X = x_i)$	0,2	0,42	0,38
$M(Y/X = x_i)$	0,5	0,686	0,768
$D(Y/X = x_i)$	0,03	0,03265	0,01163

$$MY = 0,68 \text{ (см. табл. 3)}, \quad DY = 0,0336 \text{ (см. табл. 3)}$$

Т.к. X – случайная величина, принимающая значения 2, 5 и 8 с вероятностью 0,2; 0,42 и 0,38, то такими же будут вероятности и условных математических ожиданий, и дисперсий. Т.обр., условное математическое ожидание $M(Y/X)$, так же как и условная дисперсия $D(Y/X)$ – случайные величины.

Обратим также внимание на то, что $M(Y)$, найденное в табл. 3, можно вычислить и по табл. 8 следующим образом:

$$M(Y) = M[M(Y/X)] = \sum_{i=1}^3 M(Y/X = x_i)P(X = x_i) = 0,5*0,2 + 0,686*0,42 + 0,768*0,38 = 0,68.$$

Разброс значений величины Y вокруг математического ожидания MY измеряется дисперсией $D(Y)$, или σ_Y^2 :

$$\sigma_Y^2 = D(Y) = M(Y - MY)^2. \quad (24)$$

(По табл. 3: $\sigma_Y^2 = 0,0336$.) Этот разброс может быть вызван:

- зависимостью величины Y от X (эта зависимость может быть обусловлена не только непосредственным влиянием X на Y , но и наличием случайных факторов, действующих на Y через переменную X);
- зависимостью величины Y от случайных факторов, влияющих только на Y и не влияющих на X ; эти факторы называют **остаточными**.

1) Построим показатель разброса значений величины Y , связанного с её зависимостью от фактора X .

Условное математическое ожидание $M(Y/X = x)$ является «представителем игреков», которые имеют место при $X = x$. Характеристикой разброса условных математических ожиданий $M(Y/X = x)$ относительно $M(Y)$ является дисперсия $D[M(Y/X)]$, или

$$\sigma_\phi^2 = D[M(Y/X)] = M[M(Y/X) - MY]^2 \quad (25)$$

– эта величина и будет показателем разброса значений величины Y , связанного с её зависимостью от фактора X . По таблице 16 найдём:

$$\sigma_\phi^2 = M[M(Y/X) - MY]^2 = (0,5 - 0,68)^2*0,2 + (0,686 - 0,68)^2*0,42 + (0,768 - 0,68)^2*0,38 = 0,0095.)$$

2) Теперь построим показатель разброса «игреков», связанного с влиянием остаточных факторов.

Зафиксируем какое-либо значение x величины X . Тогда причиной вариации величины Y при $X = x$ будут остаточные факторы, влияющие только на Y и не влияющие на X . Измерителем этой вариации является условная дисперсия $D(Y/X = x)$. При различных же «иксах» характеристикой разброса «игреков», вызванного влиянием на Y остаточных факторов, будет генеральное среднее из условных дисперсий, или, иначе, математическое ожидание условной дисперсии. Эту величину обозначим σ_0^2 . Имеем

$$\sigma_0^2 = M[D(Y/X)], \quad (26)$$

где при $X = x$ дисперсия $D(Y/X = x)$ вычисляется по формуле (23). (По табл. 16 найдём

$$\sigma_0^2 = M[D(Y/X)] = \sum_{i=1}^3 D(Y/X = x_i)P(X = x_i) = 0,03*0,2 + 0,03265*0,42 + 0,01163*0,38 = 0,0241.)$$

Для вычисленных дисперсий справедливо тождество

$$DY = D[M(Y/X)] + M[D(Y/X)]$$

или

$$\sigma_Y^2 = \sigma_\phi^2 + \sigma_0^2. \quad (27)$$

Степень стохастической зависимости величины Y от X измеряется **генеральным корреляционным отношением**

$$\rho_{Y/X} = + \sqrt{\frac{D[M(Y/X)]}{DY}} = + \sqrt{\frac{\sigma_\phi^2}{\sigma_Y^2}} \quad (28) = + \sqrt{1 - \frac{\sigma_0^2}{\sigma_Y^2}} = + \sqrt{1 - \frac{M[D(Y/X)]}{DY}}.$$

Квадрат корреляционного отношения

$$\rho_{Y/X}^2 = \frac{\sigma_{\varphi}^2}{\sigma_Y^2} = \frac{D[M(Y/X)]}{DY} \stackrel{(26), (25)}{=} \frac{M[M(\frac{Y}{X}) - MY]^2}{M(Y - MY)^2} \quad (29)$$

называется **генеральным коэффициентом детерминации**; он показывает, какую долю дисперсии величины Y составляет дисперсия условных математических ожиданий, или, иначе говоря, какая доля дисперсии $D(Y)$ объясняется корреляционной зависимостью Y от X . (В примере 5 $\sigma_{\varphi}^2 = 0,0095$, $\sigma_Y^2 = 0,0336$, поэтому $\rho_{Y/X}^2 = \frac{0,0095}{0,0336} = 0,28$, т.е. 28% дисперсии величины Y объясняется её корреляционной зависимостью от X ; $\rho_{Y/X} = +\sqrt{0,28} = 0,53$.)

Свойства генерального корреляционного отношения как измерителя степени корреляционной и стохастической зависимости

1. $0 \leq \rho_{Y/X} \leq 1$.

Действительно, согласно (29), $\rho_{Y/X} \geq 0$; с другой стороны, из (28) следует, что $\sigma_{\varphi}^2 \leq \sigma_Y^2$, поэтому $\rho_{Y/X} \leq 1$.

2. Условие $\rho_{Y/X} = 0$ является необходимым и достаточным для отсутствия корреляционной зависимости Y от X , т.е. для того, чтобы $M(Y/X) = \text{const}$ при любом значении x величины X .

3. Условие $\rho_{Y/X} = 1$ является необходимым и достаточным для функциональной зависимости величины Y от X .

Следствие. Чем ближе $\rho_{Y/X}$ к единице, тем в силу (28) ближе к нулю $M[D(Y/X)]$, а следовательно, и условные дисперсии $D(Y/X = x)$. Это означает, что при каждом допустимом значении x уменьшается разброс «игреков» относительно $M(Y/X = x)$. Т.обр., чем ближе $\rho_{Y/X}$ к единице, тем меньше при каждом x отличие «игреков» от постоянного числа, равного $M(Y/X = x)$, или, иначе говоря, тем выше степень стохастической зависимости Y от X . И, наоборот, чем выше степень стохастической зависимости Y от X , тем ближе $\rho_{Y/X}$ к единице.

В практических задачах наибольший интерес представляют следующие вопросы:

- существует корреляционная зависимость Y от X или нет, иначе говоря, отлично ли генеральное корреляционное отношение $\rho_{Y/X}$ от нуля или равно нулю;
- если корреляционная зависимость существует, то какой вид имеет функция регрессии (линейный, параболический или какой-либо другой).

Точно ответить на поставленные вопросы можно лишь только в том случае, когда известен закон распределения двумерной величины (X, Y) . В примере 5 этот закон задан табл. 9, в которой даны все возможные значения случайных величин X и Y и вероятности совместного появления этих значений. Обычно такими сведениями не располагают; как правило, имеются лишь наблюдавшиеся значения двумерной величины (X, Y) . Покажем как, имея наблюдавшиеся значения, ответить на поставленные выше вопросы.

Выборочное корреляционное отношение. Его значимость

Пусть имеется n наблюдений двумерной величины (X, Y) . Наблюдавшиеся «иксы» и «игреки» поместим в табл. 17, которая называется **корреляционной таблицей** и строится следующим образом:

- «иксы» группируются в вариационный ряд, число групп которого обозначим v ; если это дискретный ряд, то x_1, x_2, \dots, x_v – различающиеся между собой результаты наблюдений или варианты; если это интервальный ряд, то x_1, x_2, \dots, x_v – центры интервалов;
- «игреки» группируют в вариационный ряд, число групп которого обозначим q : y_1, y_2, \dots, y_q – это либо варианты, если ряд дискретный, либо середины интервалов, если ряд интервальный;
- подсчитывают числа m_{ji} таких наблюдавшихся пар чисел (x, y) , у которых x попадает в группу x_i , а y – в группу y_j , $i = 1, 2, \dots, v$, $j = 1, 2, \dots, q$; например, m_{12} – число пар чисел (x, y) , у которых x попало в группу x_2 , а y – в группу y_1 . Числа $m_{11}, m_{12}, \dots, m_{qv}$ называются **частотами**.

Прежде чем пояснить остальные элементы этой таблицы, сделаем следующее замечание по поводу схемы построения выборочного корреляционного отношения: от табл. 17, содержащей частоты, можно перейти к таблице частотей (табл. 18). Сравним табл.18 и табл. 9. Их различие состоит в следующем: табл.9 относилась к генеральной совокупности, поэтому в ней были указаны все мыслимые значения величин X и Y и вероятности комбинаций этих значений; табл. 18 относится к выборочной совокупности, и в ней приведены наблюдаемые значения величин X и Y и частоты, или опытные вероятности комбинаций наблюдаемых значений. Поэтому выборочное корреляционное отношение можно строить по той же схеме, что и генеральное корреляционное отношение, если заменить возможные значения величин X и Y на наблюдаемые, вероятности на частоты, математические ожидания на средние, дисперсии на выборочные дисперсии.

Однако чаще выборочное корреляционное отношение строят, используя непосредственно табл.9, а не табл.1. В табл. 9 кроме сгруппированных наблюдений и частот содержатся следующие данные:

- суммы частот по каждой строке $m_1 = \sum_{i=1}^v m_{1i}, m_2 = \sum_{i=1}^v m_{2i}, \dots, m_q = \sum_{i=1}^v m_{qi}, \quad (30)$

- суммы частот по каждому столбцу $n_1 = \sum_{j=1}^q m_{j1}, n_2 = \sum_{j=1}^q m_{j2}, \dots, n_v = \sum_{j=1}^q m_{jv}, \quad (31)$

Таблица 9.

j	i	1	2	...	v	
	$Y \quad X$	x_1	x_2	...	x_v	Σ
1	y_1	m_{11}	m_{12}	...	m_{1v}	m_1
2...	$y_{2...}$	$m_{21...}$	$m_{22...}$...	$m_{2v...}$	$m_{2...}$
q	y_q	m_{q1}	m_{q2}		m_{qv}	m_q
Σ		n_1	n_2	...	n_v	$n = \sum_{i=1}^v n_i = \sum_{j=1}^q m_j$
Групповое среднее		$\bar{Y}^{(1)}$	$\bar{Y}^{(2)}$...	$\bar{Y}^{(v)}$	
Групповая выборочная дисперсия		σ_{*1}^2	σ_{*2}^2	...	σ_{*v}^2	

Таблица 10

j	i	1	2	...	v	
	$Y \quad X$	x_1	x_2	...	x_v	
1	y_1	p_{*11}^*	p_{*12}^*	...	p_{*1v}^*	
2...	$y_{2...}$	$p_{*21...}^*$	$p_{*22...}^*$...	$p_{*2v...}^*$	$p_{*ij}^* = m_{ij}/n, i = 1, 2, \dots, v$ $j = 1, 2, \dots, q.$
q	y_q	p_{*q1}^*	p_{*q2}^*		p_{*qv}^*	

- групповые средние значения «игреков»

$$\begin{aligned} \bar{Y}^{(1)} &= (y_1 m_{11} + y_2 m_{21} + \dots + y_q m_{q1})/n_1, \\ \bar{Y}^{(2)} &= (y_1 m_{12} + y_2 m_{22} + \dots + y_q m_{q2})/n_2, \\ &\dots \end{aligned} \quad (32)$$

$$\bar{Y}^{(v)} = (y_1 m_{1v} + y_2 m_{2v} + \dots + y_q m_{qv}) / n_v.$$

Эти средние являются выборочными аналогами соответствующих условных математических ожиданий: $\bar{Y}^{(1)}$ - выборочный аналог математического ожидания величины Y при условии, что $X = x_1$, т.е. аналог величины $M(Y/X=x_1)$; $\bar{Y}^{(2)}$ - аналог $M(Y/X = x_2)$ и т.д.

Групповые выборочные дисперсии:

$$\begin{aligned} \hat{\sigma}_1^2 &= [(y_1 - \bar{Y}^{(1)})^2 m_{11} + (y_2 - \bar{Y}^{(1)})^2 m_{21} + \dots + (y_q - \bar{Y}^{(1)})^2 m_{q1}] / n_1, \\ \hat{\sigma}_2^2 &= [(y_1 - \bar{Y}^{(2)})^2 m_{12} + (y_2 - \bar{Y}^{(2)})^2 m_{22} + \dots + (y_q - \bar{Y}^{(2)})^2 m_{q2}] / n_2, \\ &\dots\dots\dots \\ \hat{\sigma}_v^2 &= [(y_1 - \bar{Y}^{(v)})^2 m_{1v} + (y_2 - \bar{Y}^{(v)})^2 m_{2v} + \dots + (y_q - \bar{Y}^{(v)})^2 m_{qv}]. \end{aligned} \quad (33)$$

Эти дисперсии являются выборочными аналогами соответствующих условных дисперсий: $\hat{\sigma}_1^2$ - выборочный аналог условной дисперсии $D(Y/X=x_1)$, $\hat{\sigma}_2^2$ - аналог $D(Y/X=x_2)$ и т.д.

Построим выборочный аналог генерального корреляционного отношения. Выборочным аналогом генеральной дисперсии $\sigma_Y^2 = DY$ является величины $\hat{\sigma}_Y^2$. Для того чтобы вычислить $\hat{\sigma}_Y^2$, найдем сначала среднее \bar{Y} по данным табл. 17. Это можно сделать по одной из следующих тождественных формул:

$$\bar{Y} = \frac{1}{n} (y_1 m_1 + y_2 m_2 + \dots + y_q m_q) = \frac{1}{n} (\bar{Y}^{(1)} n_1 + \bar{Y}^{(2)} n_2 + \dots + \bar{Y}^{(v)} n_v) \quad (34)$$

Напомним, что \bar{Y} - это выборочный аналог $M(Y)$. Теперь найдем

$$S_Y^2 = (y_1 - \bar{Y})^2 m_1 + (y_2 - \bar{Y})^2 m_2 + \dots + (y_q - \bar{Y})^2 m_q = \sum_{j=1}^q (y_j - \bar{Y})^2 m_j \quad (35)$$

Тогда

$$\hat{\sigma}_Y^2 = S_Y^2 / n \quad (36)$$

Выборочным аналогом генеральной дисперсии $\sigma_\Phi^2 = D[M(Y/X)]$ является выборочная дисперсия $\hat{\sigma}_{Y(i)}^2$ групповых средних; обозначим ее $\hat{\sigma}_\Phi^2$. Имеем

$$\hat{\sigma}_\Phi^2 = \hat{\sigma}_{Y(i)}^2 = \frac{1}{n} [(\bar{Y}^{(1)} - \bar{Y})^2 n_1 + (\bar{Y}^{(2)} - \bar{Y})^2 n_2 + \dots + (\bar{Y}^{(v)} - \bar{Y})^2 n_v] = S_\Phi^2 / n, \quad (37)$$

где

$$\begin{aligned} S_\Phi^2 &= (\bar{Y}^{(1)} - \bar{Y})^2 n_1 + (\bar{Y}^{(2)} - \bar{Y})^2 n_2 + \dots + (\bar{Y}^{(v)} - \bar{Y})^2 n_v = \\ &= \sum_{i=1}^v (\bar{Y}^{(i)} - \bar{Y})^2 n_i. \end{aligned} \quad (38)$$

Выборочным аналогом дисперсии $\sigma_o^2 = D[M(Y/X)]$ является средняя $\hat{\sigma}_i^2$ групповых выборочных дисперсий. Обозначим $\hat{\sigma}_o^2$ эту среднюю

$$\hat{\sigma}_o^2 = \overline{\hat{\sigma}_i^2} = \frac{1}{n} (\hat{\sigma}_1^2 n_1 + \hat{\sigma}_2^2 n_2 + \dots + \hat{\sigma}_v^2 n_v) = S_o^2 / n, \quad (39)$$

где $S_o^2 = \hat{\sigma}_1^2 n_1 + \hat{\sigma}_2^2 n_2 + \dots + \hat{\sigma}_v^2 n_v$.

Получаем

$$\begin{aligned} S_o^2 &= (y_1 - \bar{Y}^{(1)})^2 m_{11} + (y_2 - \bar{Y}^{(1)})^2 m_{21} + \dots + (y_q - \bar{Y}^{(1)})^2 m_{q1} + (y_1 - \bar{Y}^{(2)})^2 m_{12} + \\ &+ (y_2 - \bar{Y}^{(2)})^2 m_{22} + \dots + (y_q - \bar{Y}^{(2)})^2 m_{q2} + \dots + (y_1 - \bar{Y}^{(v)})^2 m_{1v} + (y_2 - \bar{Y}^{(v)})^2 m_{2v} + \dots \\ &+ (y_q - \bar{Y}^{(v)})^2 m_{qv} = \sum_{i=1}^v \sum_{j=1}^q (y_j - \bar{Y}^{(i)})^2 m_{ji}. \end{aligned} \quad (40)$$

Для вычисленных дисперсий справедливо тождество

$$\hat{\sigma}_Y^2 = \hat{\sigma}_{Y(i)}^2 + \hat{\sigma}_i^2, \text{ или } \hat{\sigma}_Y^2 = \hat{\sigma}_\Phi^2 + \hat{\sigma}_o^2, \quad (41)$$

аналогичное тождеству, имеющему место в генеральной совокупности.

Выборочный аналог генерального корреляционного отношения вычисляется следующим образом:

$$\hat{\rho}_{Y/X} = + \sqrt{\frac{\hat{\sigma}_{\bar{Y}(i)}^2}{\hat{\sigma}_Y^2}} = + \sqrt{\frac{\hat{\sigma}_\Phi^2}{\hat{\sigma}_Y^2}} = + \sqrt{1 - \frac{\hat{\sigma}_\Phi^2}{\hat{\sigma}_Y^2}} = + \sqrt{1 - \frac{\hat{\sigma}_1^2}{\hat{\sigma}_Y^2}} \quad (42)$$

Величина $\hat{\rho}_{Y/X}^2$ называется выборочным коэффициентом детерминации. Этот коэффициент показывает, какую долю дисперсии $\hat{\sigma}_Y^2$ составляет выборочная дисперсия групповых средних «игреков» или, иначе говоря, какая доля дисперсии $\hat{\sigma}_Y^2$ объясняется зависимостью Y от X.

Как правило, дисперсии $\hat{\sigma}_Y^2$ и $\hat{\sigma}_\Phi^2$ находят не по рассмотренным выше формулам, а по следующим, более удобным для вычислений:

$$\hat{\sigma}_Y^2 = \frac{1}{n} \sum_{j=1}^q y_j^2 m_j - (\bar{Y})^2, S_y^2 = \hat{\sigma}_Y^2 n. \quad (43)$$

$$\hat{\sigma}_i^2 = \frac{1}{n_i} \sum_{j=1}^q y_j^2 m_{ji} - (\bar{Y}^{(i)})^2, S_y^2 = \hat{\sigma}_Y^2 n. \quad (44)$$

$$\hat{\sigma}_\Phi^2 = \frac{1}{n} \sum_{i=1}^v (\bar{Y}^{(i)})^2 n_i - (\bar{Y})^2, S_\Phi^2 = \hat{\sigma}_\Phi^2 n. \quad (45)$$

$$\text{Проверим гипотезу } H_0: P_{Y/X} = 0 \quad (46)$$

Предварительно отметим, что в силу свойства 2 корреляционного отношения при выполнении гипотезы имеет место равенство условных, или групповых математических ожиданий величины Y: $M(Y/X=x_1)=M(Y/X=x_2)=\dots=M(Y/X=x_v)$.

Поэтому проверка гипотезы сводится к проверке гипотезы о равенстве групповых математических ожиданий – это задача дисперсионного анализа. Ее можно решить, если выполняются требования, применительно к нашим условиям формулирующиеся следующим образом (*):

- при каждом наблюдаемом значении x_i величины X наблюдения величины Y должны быть независимыми и проводиться в одинаковых условиях; наблюдения должны быть независимы и при различных «иксах»;
- при каждом значении x_i величина Y должна иметь нормальный закон с постоянной для различных «иксов» генеральной дисперсией (обозначим эту дисперсию σ_0^2 ;

$$\sigma_0^2 = D(Y/X = x_1) = D(Y/X = x_2) = \dots = D(Y/X = x_v)$$

Допустим, что эти требования выполняются. Тогда для проверки гипотезы (46) следует заполнить табл. 19.

В заключение заметим, что если гипотеза отвергается, то говорят, что выборочное корреляционное отношение статистически значимо. Если гипотеза не отвергается, то говорят, что выборочное корреляционное отношение незначимо.

Обратим внимание на то, что вычислить выборочное корреляционное отношение, а также проверить его значимость можно только в том случае, когда результаты наблюдений сгруппированы в таблицу типа таблицы.

Допустим, что, располагая выборочными данными, мы пришли к выводу, что корреляционная зависимость Y от X существует, т.е. при изменениях изменяются условные математические ожидания $M(Y/X=x)$. Тогда возникает вопрос: каков вид функции регрессии, т.е. функции $\phi(x) = M(Y/X=x)$?

Располагая только выборочными данными, нельзя дать точный ответ на поставленный вопрос, но высказать гипотезу о виде функции $\phi(x)$ можно; также можно провести статистическую проверку этой гипотезы, т.е. выяснить, противоречит или нет эта гипотеза имеющимся выборочным данным.

3. Виды регрессий, статистическая значимость их параметров. Автокорреляция

Исследование начинается с теории, устанавливающей связь между явлениями. Из всего круга факторов, влияющих на результативный признак, выделяются наиболее существенные факторы. После того, как было выявлено наличие взаимосвязи между

изучаемыми признаками, определяется точный вид этой зависимости с помощью регрессионного анализа.

Регрессионный анализ заключается в определении аналитического выражения (в определении функции), в котором изменение одной величины (результативного признака) обусловлено влиянием независимой величины (факторного признака). Количественно оценить данную взаимосвязь можно с помощью построения уравнения регрессии или регрессионной функции.

Базисной регрессионной моделью является модель парной (однофакторной) регрессии. Парная регрессия – уравнение связи двух переменных y и x :

$$y = f(x)$$

где y – зависимая переменная (результативный признак);

x – независимая, объясняющая переменная (факторный признак).

В зависимости от характера изменения y с изменением x различают линейные и нелинейные регрессии.

Линейная регрессия
$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Данная регрессионная функция называется полиномом первой степени и используется для описания равномерно развивающихся во времени процессов.

Наличие случайного члена ε (ошибки регрессии) связано с воздействием на зависимую переменную других неучтенных в уравнении факторов, с возможной нелинейностью модели, ошибками измерения, следовательно, появление случайной ошибки уравнения регрессии может быть обусловлено следующими объективными причинами:

1) нерепрезентативность выборки. В модель парной регрессии включается фактор, не способный полностью объяснить вариацию результативного признака, который может быть подвержен влиянию многих других факторов (пропущенных переменных) в гораздо большей степени. Например, заработная плата может зависеть, кроме квалификации, от уровня образования, стажа работы, пола и пр.;

2) существует вероятность того, что переменные, участвующие в модели, могут быть измерены с ошибкой. Например, данные по расходам семьи на питание составляются на основании записей участников опросов, которые, как предполагается, тщательно фиксируют свои ежедневные расходы. Разумеется, при этом возможны ошибки.

На основе выборочного наблюдения оценивается выборочное уравнение регрессии (линия регрессии):

$$y_x = a + b x,$$

где a , b – оценки параметров уравнения регрессии (α , β).

Аналитическая форма зависимости между изучаемой парой признаков (регрессионная функция) определяется с помощью следующих методов:

На основе теоретического и логического анализа природы изучаемых явлений, их социально-экономической сущности. Например, если изучается зависимость между доходами населения и размером вкладов населения в банки, то очевидно, что связь прямая.

Графический метод, когда характер связи оценивается визуально.

Эту зависимость можно наглядно увидеть, если построить график, отложив на оси абсцисс значения признака x , а на оси ординат – значения признака y . Нанеся на график точки, соответствующие значениям x и y , получим корреляционное поле:

а) если точки беспорядочно разбросаны по всему полю – это говорит об отсутствии зависимости между этими признаками;

б) если точки концентрируются вокруг оси, идущей от нижнего левого угла в верхний правый – то имеется прямая зависимость между признаками;

в) если точки концентрируются вокруг оси, идущей от верхнего левого угла в нижний правый – то обратная зависимость между признаками.

Если на корреляционном поле соединим точки отрезками прямой, то получим ломаную линию с некоторой тенденцией к росту. Это будет эмпирическая линия связи или эмпирическая линия регрессии. По ее виду можно судить не только о наличии, но и о форме зависимости между изучаемыми признаками.

Построение уравнения парной регрессии

Построение уравнения регрессии сводится к оценке ее параметров. Эти оценки параметров могут быть найдены различными способами. Одним из них является метод наименьших квадратов (МНК). Суть метода состоит в следующем. Каждому значению соответствует эмпирическое (наблюдаемое) значение y . Построив уравнение регрессии, например уравнение прямой линии, каждому значению будет соответствовать теоретическое (расчетное) значение y_x . Наблюдаемые значения не лежат в точности на линии регрессии, т.е. не совпадают с y_x . Разность между фактическим и расчетным значениями зависимой переменной называется остатком:

$$e = y - y_x$$

МНК позволяет получить такие оценки параметров, при которых сумма квадратов отклонений фактических значений результативного признака от теоретических y_x , т.е. сумма квадратов остатков, минимальна:

$$\sum (y - y_x)^2 \rightarrow \min$$

Для линейных уравнений и нелинейных, приводимых к линейным, решается следующая система относительно a и b :

$$an + b \sum x = \sum y$$

$$a \sum x + b \sum x^2 = \sum yx$$

где n – численность выборки.

Решив систему уравнений, получим значения a и b , что позволяет записать уравнение регрессии (регрессионное уравнение):

$$y_x = a + bx \quad \text{где } x \text{ – объясняющая (независимая) переменная;}$$

y_x – объясняемая (зависимая) переменная;

Линия регрессии проходит через точку (\bar{x}, \bar{y}) и выполняются равенства:

$$e = 0, \quad y = y_x$$

Можно воспользоваться готовыми формулами, которые вытекают из этой системы уравнений:

$$b = \frac{\text{cov}(x, y)}{\sigma_x^2} = \frac{\overline{y \cdot x} - \bar{y} \cdot \bar{x}}{\overline{x^2} - \bar{x}^2};$$

$$a = \bar{y} - b \cdot \bar{x}$$

где \bar{y} – среднее значение зависимого признака;

\bar{x} – среднее значение независимого признака;

$\overline{y \cdot x}$ – среднее арифметическое значение произведения зависимого и независимого признаков;

σ_x^2 – дисперсия независимого признака;

$\text{cov}(x, y)$ – ковариация между зависимым и независимым признаками.

Выборочной ковариацией двух переменных x , y называется средняя величина произведения отклонений этих переменных от своих средних

Параметр b при x имеет большое практическое значение и носит название коэффициента регрессии. Коэффициент регрессии показывает, на сколько единиц в среднем изменяется величина y при изменении факторного признака x на 1 единицу своего измерения.

Знак параметра b в уравнении парной регрессии указывает на направление связи:

если $b > 0$, то связь между изучаемыми показателями прямая, т.е. с увеличением факторного признака увеличивается и результативный признак, и наоборот;

если $b < 0$, то связь между изучаемыми показателями обратная, т.е. с увеличением факторного признака результативный признак y уменьшается, и наоборот.

Значение параметра a в уравнении парной регрессии в ряде случаев можно трактовать как начальное значение результативного признака y . Такая трактовка параметра a возможна только в том случае, если значение $x = 0$ имеет смысл.

После построения уравнения регрессии, наблюдаемые значения y можно представить как: $y = y_x + e$

Остатки e , как и ошибки, являются случайными величинами, однако они, в отличие от ошибок, наблюдаемы. Остаток есть та часть зависимой переменной y , которую невозможно объяснить с помощью уравнения регрессии.

На основании уравнения регрессии могут быть вычислены теоретические значения y для любых значений x .

В экономическом анализе часто используется понятие эластичности функции. Эластичность функции рассчитывается как относительное изменение y к относительному изменению x . Эластичность показывает, на сколько процентов изменяется функция при изменении независимой переменной на 1%.

Поскольку эластичность линейной функции не является постоянной величиной, а зависит от x , то обычно рассчитывается коэффициент эластичности как средний показатель эластичности.

Коэффициент эластичности показывает, на сколько процентов в среднем по совокупности изменится величина результативного признака y при изменении факторного признака на 1% от своего среднего значения:

$$\varepsilon_x = b \frac{x}{y}$$

где \bar{x} , \bar{y} – средние значения переменных x и y в выборке.

Оценка качества построенной модели регрессии

Качество модели регрессии – адекватность построенной модели исходным (наблюдаемым) данным.

Чтобы измерить тесноту связи, т.е. измерить, насколько она близка к функциональной, нужно определить дисперсию, измеряющую отклонения y от y_x и характеризующую остаточную вариацию, обусловленную прочими факторами. Они лежат в основе показателей, характеризующих качество модели регрессии.

При исследовании регрессии устанавливается однофакторная или многофакторная будет строиться модель и вид модели (линейный или нелинейный).

Обоснование вида модели состоит в выборе вида функции (некоторого аналитического выражения), с помощью которого можно будет описать изменение исследуемого показателя под воздействием факторов.

К обоснованию вида функции идут двумя путями: Теоретическим (анализируя экономическую природу x_{0i} и x_{ji} , выдвигается гипотеза о характере изменения показателя под действием фактора) И эмпирическим (закон изменения результативного показателя под действием фактора устанавливается путем анализа совокупности фактических данных по полям корреляции).

Наиболее употребительными выражениями при описании связи одного фактора и исследуемого показателя являются:

- Уравнение прямой - $x_0 = a_0 + a_1x_1$, - Уравнение параболы - $x_0 = a_0 + a_1x_1 + a_2x_1^2$, -
Уравнение гиперболы - $x_0 = a_0 + \frac{a_1}{x_1}$.

После обоснования парных взаимосвязей переходят к записи многофакторных моделей. В экономических исследованиях чаще всего применяется линейная многофакторная модель - $x_0 = a_0 + a_1x_1 + \dots + a_nx_n$.

В качестве нелинейных моделей применяются

- Мультипликативная модель - $x_0 = a_0x_1^{a_1}x_2^{a_2}x_3^{a_3} \dots$ или $x_0 = a_0a_1^{x_1}a_2^{x_2}a_3^{x_3} \dots$

Для оценки значений параметров регрессионной модели чаще всего используется Метод наименьших квадратов (МНК). Этот метод можно применить как для линейных моделей, так и для нелинейных, допускающих преобразование их к линейному виду путем замены переменных или дифференцированием.

При использовании МНК делаются определенные предпосылки относительно случайной составляющей ε . В модели $y = a + b_1 \cdot x_1 + b_2 \cdot x_2 + \varepsilon$ случайная составляющая ε представляет собой ненаблюдаемую величину. Поэтому в задачу регрессионного анализа входит не только построение самой модели, но и исследование случайных отклонений ε_i , т. е. остаточных величин.

Остатки представляют собой независимые случайные величины, и их среднее значение равно 0; они имеют одинаковую (постоянную) дисперсию и подчиняются нормальному распределению.

Статистические проверки параметров регрессии, показателей корреляции основаны на непроверяемых предпосылках распределения случайной составляющей ε_i . Связано это с тем, что оценки параметров регрессии должны отвечать определенным критериям: быть Несмещенными, состоятельными и эффективными. Эти свойства оценок, полученных по МНК, имеют чрезвычайно важное практическое значение в использовании результатов регрессии и корреляции.

Коэффициенты регрессии, найденные из системы нормальных уравнений, представляют собой выборочные оценки характеристики силы связи. Их несмещенность является желательным свойством, т. к. только в этом случае они могут иметь практическую значимость.

Несмещенность оценки означает, что математическое ожидание остатков равно нулю. Оценки считаются Эффективными, если они характеризуются наименьшей дисперсией. Поэтому несмещенность оценки должна дополняться минимальной дисперсией. Состоятельность оценок характеризует увеличение их точности с увеличением объема выработки.

Указанные критерии оценок (несмещенность, состоятельность, эффективность) обязательно учитываются при разных способах оценивания. Метод наименьших квадратов строит оценки регрессии на основе минимизации суммы квадратов остатков ($y - \hat{y}_x$).

Исследование остатков ε_i предполагают проверку наличия следующих пяти предпосылок МНК:

- случайный характер остатков; - нулевая средняя величина остатков, не зависящая от x_i ; - гомоскедастичность — дисперсия каждого отклонение ε_i одинакова для всех значений x ; - отсутствие автокорреляции остатков, т. е. значения остатков ε_i распределены независимо друг от друга; - остатки подчиняются нормальному распределению.

С целью проверки случайного характера остатков ε_i строится график зависимости остатков ε_i от теоретических значений результативного признака \hat{y} .

. Если на графике нет направленности в расположении точек ε_i , то остатки ε_i представляют собой случайные величины и МНК оправдан. Также возможны следующие случаи: если ε_i зависит от теоретического значения, то:

Вторая предпосылка МНК относительно нулевой средней величины остатков означает, что $\sum(y - \hat{y}_x) = 0$. Это выполнимо для линейных моделей и моделей, нелинейных относительно включаемых переменных. Для обеспечения несмещенности оценок коэффициентов регрессии, полученных МНК, необходимо выполнение условий независимости случайных остатков ε_i и переменных x , что исследуется в рамках соблюдения второй предпосылки МНК. С целью проверки выполнения этой предпосылки строится график зависимости случайных остатков ε от факторов, включенных в регрессию x_i . Если расположение остатков на графике не имеет направленности, то они независимы от значений x_i . Если же график показывает наличие зависимости ε_i и x_i , то модель неадекватна.

Предпосылка о нормальном распределении остатков позволяет проводить проверку параметров регрессии и корреляции с помощью критериев t и F . Вместе с тем оценки регрессии, найденные с применением МНК, обладают хорошими свойствами даже при отсутствии нормального распределения остатков, т. е. при нарушении пятой предпосылки метода наименьших квадратов.

В соответствии с третьей предпосылкой МНК требуется, чтобы дисперсия остатков была гомоскедастичной. Это означает, что для каждого значения фактора x_i остатки ε_i имеют одинаковую дисперсию. Если это условие применения МНК не соблюдается, то имеет место гетероскедастичность. Используя трехмерной изображение, рассмотрим отличие гомо - и гетероскедастичности.

Наличие гетероскедастичности будет сказываться на уменьшении эффективности оценок b_i , в частности, становится затруднительным использование формулы стандартной ошибки коэффициента регрессии, предполагающей единую дисперсию остатков для любых значений фактора.

Наличие гетероскедастичности в остатках регрессии можно проверить с помощью ранговой корреляции Спирмэна. Суть проверки заключается в том, что в случае гетероскедастичности абсолютные остатки ε_i коррелированы со значениями фактора x_i . Эту корреляцию можно измерять с помощью коэффициента ранговой корреляции Спирмэна:

$$\rho = 1 - \frac{6 \sum d^2}{n(n^2 - 1)},$$

где ρ – абсолютная разность между рангами значений x_i и $|\varepsilon_i|$.

Статистическую значимость ρ можно определить с помощью t -критерия:

$$t_{\rho} = \frac{\rho}{\sqrt{(1 - \rho^2)}} \sqrt{(n - 2)}.$$

Принято считать, что если $t_{\text{расч}} > t_{\text{табл}}$, то корреляция между ε_i и x_i статистически значима, т. е. имеет место гетероскедастичность остатков. В противном случае принимается гипотеза об отсутствии гетероскедастичности остатков.

При построении регрессионных моделей чрезвычайно важно соблюдение четвертой предпосылки МНК – отсутствие автокорреляции остатков, т. е. распределения остатков ε_i и ε_{i-1} независимы. Автокорреляция остатков означает наличие корреляции между остатками текущих и предыдущих (последующих) наблюдений. Находится коэффициент корреляции между ε_i и ε_{i-1} , и если он окажется существенно отличным от нуля, то остатки автокоррелированы и функция плотности вероятности $F(\varepsilon)$ зависит от j -ой точки наблюдения и от распределения значений остатков в других точках наблюдения.

Отсутствие автокорреляции остатков обеспечивает состоятельность и эффективность оценок коэффициентов регрессии.

До сих пор в качестве факторов рассматривались экономические переменные, принимающие количественные значения в некотором интервале. Вместе с тем может оказаться необходимым включить в модель фактор, имеющий два или более качественных уровней. Это могут быть разного рода атрибутивные признаки, такие, например, как профессия, пол, образование, климатические условия, принадлежность к определенному региону. Для того, чтобы ввести такие переменные в регрессионную модель, им должны быть присвоены те или иные цифровые метки, т. е. качественные переменные необходимо преобразовать в количественные. Такого вида сконструированные переменные в эконометрике принято называть фиктивными переменными.

Качественные признаки могут приводить к неоднородности исследуемой совокупности, что может быть учтено при моделировании двумя путями:

- регрессия строится для каждой качественно отличной группы единиц совокупности, т. е. для каждой группы в отдельности, чтобы преодолеть неоднородность единиц общей совокупности;
- общая регрессионная модель строится для совокупности в целом, учитывающей неоднородность данных. В этом случае в регрессионную модель вводятся фиктивные переменные, т. е. строится регрессионная модель с переменной структурой, отражающей неоднородность данных.

Качественный фактор может иметь только два состояния, которым будут соответствовать 1 и 0. Если же число градаций качественного признака-фактора превышает два, то в модель вводится несколько фиктивных переменных, число которых должно быть меньше числа качественных градаций. Только при соблюдении этого положения матрица исходных фиктивных переменных не будет линейно зависима и возможна оценка параметров модели.

Коэффициент регрессии при фиктивной переменной интерпретируется как среднее изменение зависимой переменной при переходе от одной категории к другой при неизменных значениях остальных параметров. На основе t -критерия Стьюдента делается вывод о значимости влияния фиктивной переменной, существенности расхождения между категориями.

Такая проверка производится с помощью статистических критериев и на их основе делается вывод о статистической надежности построенного уравнения регрессии, о пригодности модели для анализа и прогнозирования исследуемого показателя.

Для проверки надежности модели в целом используется отношение факторной

$$\frac{S_{\text{факт}}^2}{S_{\text{ост}}^2}$$

дисперсии к остаточной $S_{\text{ост}}^2$. Известно, что отношение этих дисперсий подчиняется распределению Фишера (F-распределение). Расчетное значение F-отношения сравнивается с табличным значением, которое определяется для конкретного уровня значимости α . В экономических исследованиях α принимается равным 0,05 (реже 0,01),

число степеней свободы $k_1 = p, k_2 = n - p - 1$. Если $F_{\text{расч}} > F_{\text{табл}}$, то построенная модель считается статистически надежной, а следовательно, отражает закон изменения исследуемого показателя под действием факторов.

Для проверки полноту модели используется $R^2 \cdot 100\%$. Этот показатель показывает, на сколько процентов изменится вариация результативного показателя под влиянием факторов, включенных в модель.

Проверку надежности параметров уравнения регрессии проводят с использованием Т - критерия. Расчетное значение вычисляется по формуле $t = \frac{|a_j|}{\sigma_{a_j}}$

$\sigma_{a_j} = \frac{S_{\text{ост}}}{\sigma_{x_j} \sqrt{n} \sqrt{1 - R_{j,1,2,\dots,j-1,j+1,p}^2}}$. Фактическое значение Т- критерия сравнивается с табличным и если $t_{\text{факт}} > t_{\text{табл}} (t_{\alpha,k}, \alpha = 0,05(0,01), k = n - p - 1)$, то тогда соответствующий коэффициент регрессии значим, т. е. отличен от нуля, а влияние J-го фактора следует считать сильным. Факторы, оказывающие несущественное влияние на исследуемый показатель, из модели исключают.

На этом этапе разрабатываются рекомендации об использовании результатов моделирования.

2. МЕТОДИЧЕСКИЕ УКАЗАНИЯ ПО ПРОВЕДЕНИЮ ПРАКТИЧЕСКИХ ЗАНЯТИЙ

2.1 Практическое занятие № 1 (2 часа).

Тема: «Методологическая основа научно-исследовательской работы»

2.1.1 Задание для работы:

1. Виды НИР
2. Вопросы методологии
3. Этапы НИР

2.1.2 Краткое описание проводимого занятия:

1. **Виды НИР** – доклады с презентациями, обсуждение
2. **Вопросы методологии** - доклады с презентациями, обсуждение
3. **Этапы НИР** - доклады с презентациями, обсуждение

2.1.3 Результаты и выводы:

В результате проведенного занятия студенты:

- должны ознакомиться с классификацией НИР, ее особенностями;
- усвоить основные требования к методологии научно-исследовательских работ;
- выработать навыки анализа научной проблемы и построение поэтапного плана ее решения.

2.2 Практическое занятие № 2 (2 часа).

Тема: «Математическое моделирование в инженерных исследованиях. Основные понятия и методы математической обработки экспериментальных данных»

2.2.1 Задание для работы:

1. Математическая модель: этапы построения, отличительные особенности.
2. Типовые математические модели в инженерных исследованиях.
3. Первичная обработка статистических данных. Графическое представление статистических рядов.
4. Эмпирическая функция распределения статистических рядов.
5. Числовые характеристики статистического ряда, их свойства.

2.2.2 Краткое описание проводимого занятия:

1. **Исследования в области новых материалов и технологий** – доклады с презентациями по теме, обсуждение
2. **Построение глобального информационного пространства** – доклады с презентациями по теме, обсуждение
3. **Первичная обработка статистических данных. Графическое представление статистических рядов**
4. **Эмпирическая функция распределения статистических рядов**

Пример. Записать в виде вариационного и статистического рядов выборку 5, 3, 7, 10, 5, 5, 2, 10, 7, 2, 7, 7, 4, 2, 4. Определить размах выборки.

Решение. В данном случае объем выборки $n = 15$. Упорядочим элементы выборки по величине, получим вариационный ряд 2, 2, 3, 4, 4, 5, 5, 5, 7, 7, 7, 7, 10, 10. Найдем размах выборки $\omega = 10 - 2 = 8$. Различными в заданной выборке являются элементы $z_1 = 2, z_2 = 3, z_3 = 4, z_4 = 5, z_5 = 7, z_6 = 10$; их частоты соответственно равны $n_1 = 3, n_2 = 1, n_3 = 2, n_4 = 3, n_5 = 4, n_6 = 2$. Статистический ряд исходной выборки можно записать в виде следующей таблицы:

z_i	2	3	4	5	7	10
n_i	3	1	2	3	4	2

Для контроля правильности записи находим $\sum n_i = 15$. При большом объеме выборки ее элементы рекомендуется объединять в группы (разряды), представляя результаты опытов в виде *группированного статистического ряда*. В этом случае интервал, содержащий все элементы выборки, разбивается на k непересекающихся интервалов. Вычисления упрощаются, если эти интервалы имеют одинаковую длину $b \approx \frac{\omega}{k}$. В дальнейшем рассматривается именно этот случай. После того как частичные интервалы выбраны, определяют частоты - количество n_i элементов выборки, попавших в i -й интервал (элемент, совпадающий с верхней границей интервала, относится к следующему интервалу). Получающийся статистический ряд в верхней строке содержит середины z_i интервалов группировки, а в нижней — частоты n_i ($i = 1, 2, \dots, k$).

Наряду с частотами одновременно подсчитываются также накопленные частоты $\sum_{j=1}^i n_j$, относительные частоты n_i/n и *накопленные относительные частоты* $\sum_{j=1}^i n_j/n$, $i = 1, 2, \dots, k$. Полученные результаты сводятся в таблицу, называемую *таблицей частот группированной выборки*.

Пример. Представить выборку 55 наблюдений в виде таблицы частот, разбив имеющиеся данные выборки на семь интервалов группировки. Выборка:

20,3	15,4	17,2	19,2	23,3	18,1	21,9
15,3	16,8	13,2	20,4	16,5	19,7	20,5
14,3	20,1	16,8	14,7	20,8	19,5	15,3
19,3	17,8	16,2	15,7	22,8	21,9	12,5
0,1	21,1	18,3	14,7	14,5	18,1	18,4
13,9	19,1	18,5	20,2	23,8	16,7	20,4
19,5	17,2	19,6	17,8	21,3	17,5	19,4
17,8	13,5	17,8	11,8	18,6	19,1	

В данном случае размах выборки $\omega = 23,8 - 0,1 = 23,7$; тогда длина интервала группировки будет $b = 23,7/7 \approx 3,4$. В качестве первого интервала возьмем интервал 10 - 12. Результаты группировки сведем в таблицу 1

Таблица 1

Номер интервала i	Границы интервала	Середина интервала z_i	Частота n_i	Накопленная частота $\sum_{j=1}^i n_j$	Относительная частота n_i/n	Накопленная относительная частота $\sum_{j=1}^i n_j/n$
1	10-12	11	2	2	0,0364	0,0364
2	12-14	13	4	6	0,0727	0,1091
3	14-16	15	8	14	0,1455	0,2546
4	16-18	17	12	26	0,2182	0,4728
5	18-20	19	16	42	0,2909	0,7637
6	20-22	21	10	52	0,1818	0,9455
7	22-24	23	3	55	0,0545	1,0000

Пример. Построить гистограмму и полигон частот, а также график эмпирической функции распределения группированной выборки.

Решение. По результатам группировки (см. таблицу 1.) строим гистограмму частот (рис. 1). Соединяя отрезками ломаной середины верхних оснований прямоугольников, из которых состоит полученная гистограмма, получаем соответствующий полигон частот (рис. 2). Так как середина первого интервала группировки $z_1 = 11$, то $F_n^*(x) = 0$ при $x \leq 11$. Рассуждая аналогично, находим, что $F_n^*(x) = 1$ при $x > 23$. На полуинтервале $(11, 23]$ эмпирическую функцию распределения строим по данным третьего и последнего столбцов таблицы 1.

$F_n^*(x)$ имеет скачки в точках, соответствующих серединам интервалов группировки. В результате получаем график $F_n^*(x)$, изображенный на рис. 3.

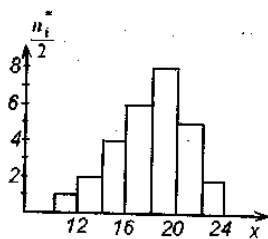


Рис.1

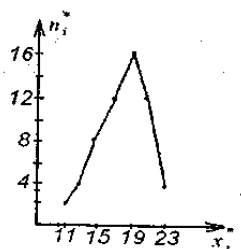


Рис.2

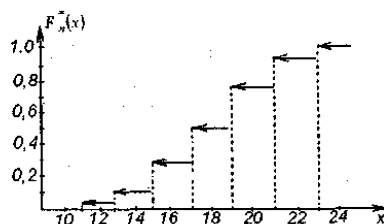


Рис.3

5. Числовые характеристики статистического ряда, их свойства.

Пример. Определить среднее, моду и медиану для выборки 5, 6, 8, 2, 3, 1, 1, 4.

Решение. Представим данные в виде вариационного ряда: 1, 1, 2, 3, 4, 5, 6, 8. Выборочное

среднее $\bar{x} = \frac{1}{8}(1+1+2+3+4+5+6+8) = 3,75$. Все элементы входят в выборку по одному

разу, кроме 1, следовательно, мода $\tilde{d}_x = 1$. Так как $n = 8$, то медиана $\tilde{h}_x = \frac{1}{2}(3+4) = 3,5$.

Итак, $\bar{x} = 3,75$, $\tilde{d}_x = 1$, $\tilde{h}_x = 3,5$.

Для упрощения вычислений выборочных среднего и дисперсии группированной выборки,

эту выборку преобразуют так: $u_i = \frac{1}{b}(z_i - d_x^*)$, $i = 1, 2, \dots, k$, где d_x^* - выборочная мода, а b -

длина интервала группировки. Эти соотношения показывают, что в выборку z_1, z_2, \dots, z_n

внесена систематическая ошибка d_x^* , а результат подвергнут преобразованию масштаба с

коэффициентом $k = 1/b$. Полученный в результате набор чисел u_1, u_2, \dots, u_n можно

рассматривать как выборку из генеральной совокупности $U = \frac{1}{b}(x - d_x^*)$. Тогда

выборочные среднее \bar{x} и дисперсия D_x^* исходных данных связаны со средним \bar{u} и

дисперсией D_u^* преобразованных данных следующими соотношениями: $\bar{x} = b\bar{u} + d_x^*$,

$D_x^* = b^2 D_u^*$.

Пример. Вычислить среднее и дисперсию группированной выборки

Границы интервалов	134-138	138-142	142-146	146-150	150-154	154-158
Частоты	1	3	15	18	14	2

Решение. Длина интервала группировки $b = 4$, значение середины интервала,

встречающегося с наибольшей частотой $d_x^* = 148$. Преобразование последовательности

середин интервалов выполняется по формуле:

$$u_i = \frac{z_i - 148}{4}, \text{ где } i = 1, 2, \dots, 6.$$

Таблица 2

i	z_i	u_i	n_i	$n_i u_i$	$n_i u_i^2$	$n_i (u_i + 1)^2$
1	136	-3	1	-3	9	4
2	140	-2	3	-6	12	3
3	144	-1	15	-15	15	0
4	148	0	18	0	0	18
5	152	1	14	14	14	56
6	156	2	2	4	8	18
Σ	-	-	53	-6	58	99

Вычисления сведены в таблицу 2. Последний столбец этой таблицы служит для контроля вычислений при помощи тождества $\sum n_i (u_i + 1)^2 = \sum n_i u_i^2 + 2 \sum n_i u_i + \sum n_i$. Выполняя вычисления, получим $58 + 2 \cdot (-6) + 23 = 99$.

Полученный результат показывает, что вычисления выполнены правильно. По формулам, данным выше, находим средние значения U

$$\bar{u} = \frac{-6}{53} \approx -0,113, \quad D_U^* = \frac{58 - (-6)^2 / 53}{53} \approx 1,108. \text{ Далее находим средние данной выборки:}$$

$$\bar{x} \approx (-0,113) \cdot 4 + 148 \approx 147,548, \quad D_x^* \approx 4^2 \cdot 1,103 \approx 17,728.$$

Пример. Вычислить среднее, дисперсию, коэффициенты асимметрии и эксцесса для следующей группированной выборки:

Границы интервалов	10-12	12-14	14-16	16-18	18-20	20-22	22-24
Частоты	2	4	8	12	16	10	3

Длина интервала группировки $b = 2$. Значение z_i , встречающееся с наибольшей частотой, $d_x^* = 19$. Поэтому преобразование имеет вид $u_i = \frac{z_i - 19}{2}$, где $i = 1, 2, \dots, 7$.

Все вычисления оформим в виде таблицы 3.

Таблица 3

i	z_i	u_i	n_i	$n_i u_i$	u_i^2	$n_i u_i^2$	u_i^3	$n_i u_i^3$	u_i^4	$n_i u_i^4$	$u_i + 1$	$(u_i + 1)^4$	$n_i (u_i + 1)^4$
1	11	-4	2	-8	16	32	-64	-128	256	512	-3	81	162
2	13	-3	4	-12	9	36	-27	-108	81	324	-2	16	64
3	15	-2	8	-16	4	32	-8	-64	16	128	-1	1	8
4	17	-1	12	-12	1	12	-1	-12	1	12	0	0	0
5	19	0	16	0	0	0	0	0	0	0	1	1	16
6	21	1	10	10	1	10	1	10	1	10	2	16	160
7	23	2	3	6	4	12	8	24	16	48	3	81	243
Σ	-	-	55	-32	-	134	-	-278	-	1034	-	-	653
$\Sigma \frac{n_i u_i^4}{n}$			-	-0,582	-	2,436	-	-5,054	-	18,8	-	-	-

Контроль вычислений будет $54 + 4 \cdot (-32) + 6 \cdot 134 + 4 \cdot (-278) + 1034 = 653$.

Находим искомые характеристики выборочного распределения:

$$\bar{x} = 19 + 2 \cdot \frac{-32}{55} \approx 17,8, \quad D_U^* = \frac{134 - (-32)^2 / 55}{55} \approx 2,10, \quad D_X^* = 2^2 \cdot 2,10 = 8,40,$$

$$a_X^* = \frac{1}{2,10^{3/2}} \left[-5,054 - 3 \cdot (-0,582) \cdot 2,436 + 2(-0,582)^3 \right] \approx -0,393,$$

2.2.3 Результаты и выводы:

В результате проведенного занятия студенты:

- рассмотреть типовые математические модели, применяемые в инженерных приложениях и условия их использования;
- усвоить основные понятия, связанные с математическим моделированием;
- выработать навыки анализа этапов построения математических моделей;
- должны ознакомиться с основными понятиями математической статистики, ее предметом и задачами;
- усвоить алгоритмы первичного статистического анализа экспериментальных данных;
- выработать навыки нахождения числовых характеристик статистического ряда.

2.3 Практическое занятие № 3 (2 часа).

Тема: «Основы корреляционно-регрессионного анализа»

2.3.1 Задание для работы:

1. Многомерные СВ, законы их распределения, условные числовые характеристики
2. Функция регрессии, коэффициент детерминации, корреляции, ковариация

2.3.2 Краткое описание проводимого занятия:

1. Многомерные СВ, законы их распределения, условные числовые характеристики

Пример 1. Найти выборочные средние, дисперсии и коэффициент корреляции для выборки, приведенной в таблице. Построить диаграмму рассеивания.

Решение. Вычисление указанных выборочных характеристик удобно выполнять в следующей последовательности. Сначала вычисляют суммы $\sum x_i$, $\sum y_i$, $\sum x_i^2$, $\sum y_i^2$, $\sum x_i y_i$, $\sum (x_i + y_i)^2$. Для контроля правильности вычислений используется тождество $\sum (x_i + y_i)^2 = \sum x_i^2 + 2\sum x_i y_i + \sum y_i^2$.

Таблица 1

x	y	x	y	x	y	x	y	x	y
8,35	3,50	10,50	6,00	11,35	9,50	12,15	6,00	12,85	9,50
8,74	1,49	10,75	2,50	11,50	6,00	12,25	8,05	13,15	9,02
9,25	6,40	10,76	5,74	11,50	9,00	12,35	5,01	13,25	6,49

9,50	4,50	11,00	8,50	11,62	8,50	12,50	7,03	13,26	10,50
9,75	5,00	11,00	5,26	11,75	10,00	12,76	7,53	13,40	7,51
10,24	7,00	11,25	8,00	12,00	9,00	12,85	6,01	13,50	10,00
13,65	9,50	14,50	10,00	13,75	8,51	14,75	12,00	14,00	11,00
15,25	12,50	14,23	8,40	16,00	11,50	14,26	10,00	16,00	13,00
14,51	9,50	16,25	12,00						

Объем выборки $n = 42$.

Выборочные средние отсюда находятся по формулам

$$\bar{x} = \alpha_{1,0}^* = \frac{1}{n} \sum x_i, \quad \bar{y} = \alpha_{0,1}^* = \frac{1}{n} \sum y_i$$

Затем вычисляются суммы квадратов отклонений от среднего и произведений отклонений

$$\text{от средних: } Q_x = \sum (x_i - \bar{x})^2 = \sum x_i^2 - \frac{(\sum x_i)^2}{n}, \quad Q_y = \sum (y_i - \bar{y})^2 = \sum y_i^2 - \frac{(\sum y_i)^2}{n},$$

$$Q_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n}$$

$$\text{Отсюда } D_X^* = \frac{1}{n} Q_x, \quad D_Y^* = \frac{1}{n} Q_y, \quad r = \frac{\mu_{1,1}^*}{\sqrt{D_X^* D_Y^*}} = \frac{Q_{xy}}{\sqrt{Q_x Q_y}}.$$

Предварительно вычислим

$$\sum x_i = 522,23, \quad \sum y_i = 336,41, \quad \sum x_i^2 = 6652,25, \quad \sum y_i^2 = 2987,80, \quad \sum x_i y_i = 4358,626.$$

Тогда найдем $\bar{x} = 12,434, \quad \bar{y} = 8,011$.

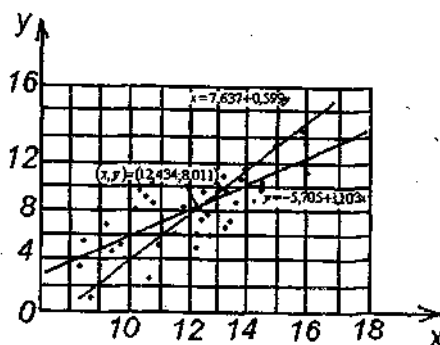
$$\text{Далее находим } Q_x = 6652,25 - \frac{522,23^2}{42} \approx 158,8182, \quad Q_y = 2987,805 - \frac{336,41^2}{42} \approx 292,5958,$$

$$Q_{xy} = 4358,626 - \frac{522,23 \cdot 336,41}{42} \approx 175,1912.$$

Окончательно, получаем

$$D_X^* = \frac{158,8182}{42} \approx 3,7814, \quad D_Y^* = \frac{292,5958}{42} \approx 6,9666, \quad r = \frac{175,1912}{\sqrt{158,8182 \cdot 292,5958}} \approx 0,813.$$

Диаграмма рассеивания приведена на рис. 1.



Выборочная линейная регрессия Y на X по выборке $(x_i, y_i), i=1,2,\dots,n$, определяется

$$\text{уравнением } y = \beta_0^* + \beta_1^* x = \bar{y} + r \frac{D_Y^*}{D_X^*} (x - \bar{x}).$$

Коэффициенты β_0^* и β_1^* называются *выборочными коэффициентами регрессии*. Они

вычисляются по формулам: $\beta_1 = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{n \sum x_i^2 - (\sum x_i)^2} = \frac{Q_{xy}}{Q_x}$, $\beta_0 = \bar{y} - \beta_1 \bar{x}$.

Аналогично определяется выборочная линейная регрессия X на Y :

$x = \beta_0^* + \beta_1^* y = \bar{x} + r \frac{D_x^*}{D_y^*} (y - \bar{y})$ коэффициенты β_0^* и β_1^* которой находятся по формулам

$$\beta_1^* = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{n \sum y_i^2 - (\sum y_i)^2} = \frac{Q_{xy}}{Q_y}, \quad \beta_0^* = \bar{x} - \beta_1^* \bar{y}.$$

Для контроля правильности расчетов используют соотношение $\sqrt{\beta_1^* \beta_1} = |r|$.

Прямые $y = \beta_0^* + \beta_1^* x$, $x = \beta_0^* + \beta_1^* y$ пересекаются в точке с координатами (\bar{x}, \bar{y}) .

Пример 2. Вычислить выборочные коэффициенты линейной регрессии X на Y и Y на X по выборке из предыдущего примера. Нанести прямые регрессии на диаграмму рассеивания.

Решение. Воспользуемся результатами вычислений в предыдущем примере. По формулам находим

$$\beta_1^* = \frac{175,1912}{158,8182} \approx 1,103, \quad \beta_0^* = 8,011 - 1,103 \cdot 12,434 \approx -5,705.$$

Таким образом, прямая регрессии Y на X имеет уравнение $y = -5,705 + 1,103x$.

Аналогично находим $\beta_1^* = \frac{175,1912}{292,5958} \approx 0,599$, $\beta_0^* = 12,434 - 0,599 \cdot 8,011 \approx 7,637$.

Отсюда прямая регрессии X на Y имеет уравнение $x = 7,637 + 0,599y$.

Проверка показывает $\sqrt{1,103 \cdot 0,599} \approx 0,813$, что полученный результат совпадает со значением r , вычисленным в примере (*). Прямые регрессии нанесены на диаграмму рассеивания на рис.1.

Пример 3. Используя группировку выборки, заданной таблицей в примере (*), вычислить выборочные средние, дисперсии, коэффициент корреляции, а также выборочные коэффициенты линейной регрессии X на Y и Y на X .

Решение. Выберем $b_x=1$, $b_y=2$. Прямоугольная сетка, соответствующая этим значениям, нанесена на диаграмму рассеивания (рис. 1). Непосредственно по диаграмме строим корреляционную таблицу (таблица 2). Находим $d_x^* = 11,5$, $d_y^* = 9$ и вычисляем значения u_i

и v_j по формулам $u_i = \frac{\hat{x}_i - 11,5}{1}$, $i=1,2,\dots,9$, $v_j = \frac{\hat{y}_j - 9}{2}$, $j=1,2,\dots,7$.

Вычисляем следующие суммы:

$$\sum n_i u_i = 43, \quad \sum n_j v_j = -15, \quad \sum n_i u_i^2 = 215, \quad \sum n_j v_j^2 = 87, \quad \sum_{i=1}^9 \sum_{j=1}^7 n_{ij} u_i v_j = 80. \quad \text{По формулам}$$

находим $Q_u = 215 - \frac{43^2}{42} \approx 170,976$, $Q_v = 87 - \frac{(-15)^2}{42} \approx 81,643$, $Q_{uv} = 80 - \frac{43 \cdot (-15)}{42} \approx 95,357$.

Далее получаем $\bar{x} = 1 \frac{43}{42} + 11,5 \approx 12,52$, $\bar{y} = 2 \frac{(-15)}{42} + 9 \approx 8,28$, $D_x^* = 1^2 \frac{170,976}{42} \approx 4,071$,

$$D_y^* = 2^2 \frac{81,643}{42} \approx 7,775, \quad r = \frac{95,357}{\sqrt{170,976 \cdot 81,643}} \approx 0,807. \quad \text{Находим выборочные коэффициенты}$$

регрессии: $\beta_1^* = \frac{2}{1} \frac{95,357}{170,976} \approx 1,12$, $\beta_1^* = \frac{1}{2} \frac{95,357}{81,643} \approx 0,58$, $\beta_0^* = 8,28 - 1,12 \cdot 12,52 \approx -5,74$,

$$\beta_0^* = 12,52 - 0,58 \cdot 8,28 \approx 7,72 .$$

Окончательно получим, что уравнение линейной регрессии Y на X имеет вид $y = -5,74 + 1,12x$, а уравнение линейной регрессии X на Y имеет вид $x = 7,72 + 0,58y$.

Расхождение полученных результатов с результатами выше рассмотренных примеров обусловлено группировкой.

Таблица 2. Корреляционная таблица для диаграммы рассеивания

Границы и середины интервалов для у	v_j	Границы и середины интервалов для x									n_j	$n_j v_j$	$n_j v_j^2$
		8-9 8,5	9- 10 9,5	10- 11 10,5	11- 12 11,5	12- 13 12,5	13- 14 13,5	14- 15 14,5	15- 16 15,5	16- 17 16, 5			
		u_i											
		-3	-2	-1	0	1	2	3	4	5			
0-2 1	-4	1	0	0	0	0	0	0	0	0	1	-4	16
2-4 3	-3	1	0	1	0	0	0	0	0	0	2	-6	18
4-6 5	-2	0	2	1	1	1	0	0	0	0	5	-10	20
6-8 7	-1	0	1	2	1	4	2	0	0	0	10	-10	10
8-10 9	0	0	0	0	5	3	3	2	0	0	13	0	0
10-12 11	1	0	0	0	1	0	2	3	0	1	7	7	7
12-14 13	2	0	0	0	0	0	0	1	1	2	4	8	16
n_i		2	3	4	8	8	7	6	1	3	$\Sigma=42$	$\Sigma=-15$	$\Sigma=87$
$n_i v_i$		-6	-6	-4	0	8	14	18	4	15	$\Sigma=43$		
$n_i v_i^2$		18	12	4	0	8	28	54	16	75	$\Sigma=215$		

2. Функция регрессии, коэффициент детерминации, корреляции, ковариация

Допустим, что в результате лечения 12 больных с артериальной гипертензией в результате суточного мониторингирования систолического артериального давления (САД) до лечения и после месячного лечения были получены следующие результаты:

№	САД до (x_i)	САД после (y_i)	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
			-9,6	-7,5	
			-19,6	-17,5	
			-14,6	-12,5	182,5
			-4,6	7,5	-34,5
			0,4	12,5	
			5,4	12,5	67,5
			-9,6	-7,5	

			10,4	7,5	
			15,4	12,5	192,5
			0,4	-7,5	-3
			5,4	-2,5	-13,5
			20,4	2,5	
	$\bar{x} = 169,6$	$\bar{y} = 137,5$			$\Sigma = 1012,5$
	$\sigma_x = 12,1$	$\sigma_y = 10,6$			

$$r_{x,y} = \frac{1}{12-1} \cdot \frac{1012,5}{12,1 \cdot 10,6} = 0,718$$

Итак, коэффициент корреляции получился равным 0,718.

Определим, достоверно ли он отличается от нуля. Для этого используем Таблицу 10 приложения. У нас 12 пар измерений, поэтому входим в Таблицу по 12 строке. На пересечении 12 строки и столбца $P=0,05$ стоит число 0,576. Полученный коэффициент корреляции (0,718) больше этого числа.

Следовательно, на этом уровне коэффициент корреляции достоверно отличается от нуля, то есть связь есть. На пересечении этой же строки и столбца $P=0,01$ стоит число 0,708. Поскольку коэффициент корреляции больше и этого числа, следовательно, мы можем говорить, что связь существует и на этом более значимом уровне. Итак, ответ на первый вопрос таков: существование связи высоко достоверно. Далее, поскольку получено положительное значение коэффициента корреляции, мы заключаем, что связь прямая. Используя Таблицу 2 данного раздела, мы приходим к заключению, что связь сильная. **Найдем коэффициент детерминации:**

$d = r^2 \times 100 = 0,718^2 \times 100 = 0,516 \times 100 = 51,6 (\%)$ Таким образом, систолическое артериальное давление после лечения на 51,6 % определяется систолическим артериальным давлением до лечения, а на 48,4 % другими факторами.

Формы проявления взаимосвязей явлений и процессов весьма разнообразны. Из них в самом общем виде выделяют *функциональную* (полную) и *корреляционную* (неполную) связи. Математически ковариация представляет собой меру линейной зависимости двух случайных величин. Коэффициент корреляции - это математическая мера корреляции двух величин. Коэффициенты корреляции могут быть положительными и отрицательными.

Иногда показателям тесноты связи можно дать качественную оценку (шкала Чеддока):

Количественная мера связи	Качественная тесноты характеристика силы связи
0,1 - 0,3	Слабая
0,3 - 0,5	Умеренная
0,5 - 0,7	Заметная
0,7 - 0,9	Высокая
0,9 - 0,99	Весьма высокая

2.3.3 Результаты и выводы:

В результате проведенного занятия студенты должны:

- ознакомиться с основными понятиями многомерного статистического анализа, теории корреляции, классификацией регрессий;
- усвоить алгоритмы нахождения условных законов и числовых характеристик многомерных случайных величин, вычисления коэффициента корреляции, детерминации, ковариации;
- выработать навыки нахождения уравнения регрессии, проверки его параметров на статистическую значимость.