

Раздел 2.

Математическая статистика

Лектор: старший преподаватель кафедры математики
Константиновская Наталья Валерьевна



Тема. Основы выборочного метода

План:

1. Основные понятия математической статистики.
2. Числовые характеристики выборки.
3. Оценка параметров генеральной совокупности по ее выборке.

1. Основные понятия математической статистики

Математическая статистика – это раздел математики, изучающий способы сбора статистической информации и методы ее обработки.

В математической статистике выделяют два основных направления исследований:

1. Оценка параметров генеральной совокупности.
2. Проверка статистических гипотез.

Генеральная совокупность – это множество всех изучаемых объектов.

Выборочная совокупность (выборка) – это часть генеральной совокупности, выбранная некоторым (случайным) образом.

Объемом совокупности (выборочной или генеральной) называют число объектов этой совокупности.

Например, из десяти тысяч студентов отобрано для обследования 100 человек.

Объем генеральной совокупности $N=10000$;

объем выборки $n=100$.

Выборка должна быть репрезентативной, то есть давать правильное представление о пропорциях генеральной совокупности.

Выборка будет репрезентативной, если ее осуществить случайно, все объекты имеют одинаковую вероятность попасть в выборку.

Каждый элемент выборки x_i
называется **вариантой**.

Число наблюдений варианты n_i
называется **частотой встречаемости**
(частотой).

Относительная частота – это
отношение частоты к объему выборки

$$W_i = \frac{n_i}{n}$$

Статистическим распределением выборки называют перечень вариантов и соответствующих им частот или относительных частот.

x_i	x_1	x_2	...	x_k
n_i	n_1	n_2	...	n_k

Гистограмма частот – это ступенчатая фигура, состоящая из смежных прямоугольников, построенных на одной прямой, основаниями которых служат частичные интервалы длины h , а высоты равны отношению n_i к h .

Полигон частот – ломаная линия, отрезки которой соединяют точки с координатами $(x_i; n_i)$.

2. Числовые характеристики выборки

Характеристики положения

Мода (M_0) – это такое значение варианты, что предшествующее и следующее за ним значения имеют меньшие частоты встречаемости.

Для одномодальных распределений мода – это наиболее часто встречающаяся варианта в данной совокупности.

Медиана (M_E) - это значение признака, относительно которого ряд распределения делится на 2 равные по объему части.

Выборочная средняя – это среднее арифметическое значение вариант статистического ряда

$$\bar{x}_B = \frac{1}{n} \sum_{i=1}^k x_i n_i$$

Характеристики рассеяния вариант вокруг своего среднего

Выборочная дисперсия – среднее арифметическое квадратов отклонения вариант от их среднего значения

$$D_B = \frac{1}{n} \sum_{i=1}^k (x_i - \bar{x}_B)^2 \cdot n_i$$

Выборочная дисперсия может быть подсчитана по формуле

$$D_B = \frac{1}{n} \sum_{i=1}^n x_i^2 \cdot n_i - (\bar{x}_B)^2$$

Среднее квадратическое отклонение – это квадратный корень из выборочной дисперсии

$$\sigma_B = \sqrt{D_B}$$

3. Оценка параметров генеральной совокупности по ее выборке

Числовые значения, характеризующие генеральную совокупность, называются параметрами.

Одна из задач математической статистики – определение параметров большого массива по исследованию его части.

Статистическое оценивание может выполняться двумя способами:

1). Точечная оценка – оценка, которая дается для некоторой определенной точки.

2). Интервальная оценка – по данным выборки оценивается интервал, в котором лежит истинное значение с заданной вероятностью.

Несмещенной оценкой генеральной средней (математического ожидания) служит выборочная средняя

$$\bar{x}_{ГЕН} \cong \bar{x}_B$$

Выборочная дисперсия не обладает свойством несмещенности.

Это смещенная оценка генеральной дисперсии .

На практике используют исправленную выборочную дисперсию, которая является несмещенной оценкой дисперсии генеральной совокупности:

$$S^2 = \frac{n \cdot D_B}{n - 1}$$

$$D_{ГЕН} \cong S^2$$

Кроме того, в расчетах используют S - исправленное среднее квадратическое отклонение, называемое стандартным отклонением

$$\sigma_{ГЕН} \cong S$$

Пример. Найти точечные оценки генеральной совокупности по данному статистическому распределению выборки.

x_i	3-5	5-7	7-9	9-11
n_i	5	10	20	15

Решение. Дан интервальный ряд распределения. Составим дискретный ряд, находя середины интервалов

x_i	4	6	8	10
n_i	5	10	20	15

$$\bar{x}_B = \frac{4 \cdot 5 + 6 \cdot 10 + 8 \cdot 20 + 10 \cdot 15}{50} = 7,8$$

$$\Rightarrow \bar{x}_T \cong 7,8$$

$$D_B = \frac{4^2 \cdot 5 + 6^2 \cdot 10 + 8^2 \cdot 20 + 10^2 \cdot 15}{50} - (7,8)^2 = 3,56$$

$$S^2 = \frac{50 \cdot 3,56}{50 - 1} = 3,63 \quad \Rightarrow D_{\Gamma} \cong 3,63$$

$$\sigma_{\Gamma} \cong \sqrt{3,63} = 1,906$$

Тема. Проверка статистических гипотез

План:

1. Основные понятия теории статистических гипотез.
2. Общая постановка задачи проверки гипотез.
3. Проверка гипотез относительно средних (критерий Стьюдента).
4. Проверка гипотез о законах распределения.

1. Основные понятия теории статистических гипотез

Статистическая гипотеза – это любое предположение о виде неизвестного распределения или о параметрах известных распределений.

Статистическая гипотеза – это всякое высказывание о генеральной совокупности, проверяемое по выборке.

Процедура сопоставления
высказанного предположения
(гипотезы) с выборочными
данными называется проверкой
гипотез.

Гипотезы будем обозначать буквой H с индексами. Будем предполагать, что у нас имеется 2 непересекающиеся гипотезы H_0 и H_1 .

H_0 – нулевая гипотеза (или основная).

H_1 – альтернативная или конкурирующая гипотеза.

Выдвинутая гипотеза может быть правильной или неправильной, поэтому возникает необходимость ее проверки.

Задача проверки статистических гипотез состоит в том, чтоб на основе выборки $x_1, x_2, x_3, \dots, x_n$

принять (т. е. считать справедливой) либо нулевую гипотезу , либо конкурирующую гипотезу .

При проверке гипотезы может быть принято неправильное решение, то есть могут быть допущены ошибки двух родов:

Ошибка первого рода состоит в том, что отвергается нулевая гипотеза H_0 , когда на самом деле она верна.

Ошибка второго рода состоит в том, что отвергается альтернативная гипотеза H_1 , когда на самом деле она верна.

Рассматриваемые случаи наглядно иллюстрирует следующая таблица.

Гипотеза H_0	Отвергается	Принимается
верна	ошибка 1-го рода	правильное решение
неверна	правильное решение	ошибка 2-го рода

Вероятность ошибки первого рода называется уровнем значимости критерия.

Для проверки принятой гипотезы используют статистический критерий – это правило, позволяющее, основываясь только на выборке $x_1, x_2, x_3, \dots, x_n$, принять либо отвергнуть нулевую гипотезу .

Различают два вида критериев: параметрические и непараметрические.

Параметрические критерии представляют собой функции параметров данной совокупности и используются, если совокупности, из которых взяты выборки, подчиняются нормальному закону распределения.

Непараметрические критерии применяются, если нет подчинения распределения нормальному закону.

2. Общая постановка задачи проверки гипотез

1. Формулируют (выдвигают) нулевую гипотезу об отсутствии различий между группами, об отсутствии существенного отличия фактического распределения от некоторого заданного, например, нормального, экспоненциального и др.

Сущность нулевой гипотезы :
разница между сравниваемыми
генеральными параметрами равна
нулю, и различия, наблюдаемые
между выборочными
характеристиками, носят случайный
характер, то есть эти выборки
принадлежат одной генеральной
совокупности.

2. Формулируют противоположную нулевой альтернативную гипотезу .

3. Задают уровень значимости α .

Уровень значимости - это вероятность ошибки отвергнуть нулевую гипотезу , если на самом деле эта гипотеза верна.

При $\alpha \leq 0,05$ ошибка возможна в 5% случаев.

4. Для проверки выдвинутой гипотезы используют критерии.

Критерий – это случайная величина K , которая служит для проверки H_0 . Эти функции распределения известны и табулированы.

Критерий зависит от двух параметров: от числа степеней свободы и от уровня значимости. Фактическую величину критерия получают по данным наблюдения $K_{НАБЛ}$.

5. По таблице определяют критическое значение, превышение которого при справедливости гипотезы маловероятно $K_{КРИТ}(\alpha, f)$

6. Сравнивают $K_{НАБЛ}$ и $K_{КРИТ}(\alpha, f)$.

Если $K_{НАБЛ} > K_{КРИТ}(\alpha, f)$, то отвергают H_0 и принимают H_1 .

Если $K_{НАБЛ} < K_{КРИТ}(\alpha, f)$, то отвергают H_1 и принимают H_0 .

7. Вывод: различие статистически значимо (0,05) или незначимо.

3. Проверка гипотез относительно средних

Сравнивают друг с другом две независимые выборки объемов n_1 и n_2 , взятые из нормально распределенных совокупностей с параметрами $M(X_1)$ и $M(X_2)$. Дополнительно предполагаем, что неизвестные генеральные дисперсии равны между собой. По этим выборкам найдены соответствующие выборочные средние \bar{x}_1 и \bar{x}_2

и исправленные дисперсии S_1^2 и S_2^2 . Уровень значимости задан.

1. Нулевая гипотеза $H_0: M(X_1) = M(X_2)$;
2. Альтернативная гипотеза $H_1: M(X_1) \neq M(X_2)$
3. $\alpha \leq 0,05$
4. Для проверки нулевой гипотезы в этом случае можно использовать критерий Стьюдента сравнения средних.

Величину критерия находим по формуле:

$$t_{\text{НАБЛ}} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}} \cdot \sqrt{\frac{n_1 \cdot n_2 (n_1 + n_2 - 2)}{n_1 + n_2}}$$

Доказано, что величина $t_{НАБЛ}$ при справедливости нулевой гипотезы имеет

t – распределение Стьюдента с

$$f = n_1 + n_2 - 2$$

степенями свободы.

5. По таблице находим $t_{КРИТ}(\alpha, f = n_1 + n_2 - 2)$

6. Сравниваем $t_{КРИТ}$ и $t_{НАБЛ}$.

Если $|t_{НАБЛ}| < t_{КРИТ}(\alpha, f) \Rightarrow H_0$

Если $|t_{НАБЛ}| > t_{КРИТ}(\alpha, f) \Rightarrow H_1$ различие
достоверно

Пример.

По двум независимым малым выборкам объемов $n_1=5$ и $n_2=6$, извлеченным из нормальных генеральных совокупностей X_1 и X_2 , вычислены выборочные средние:

$$\bar{x}_1 = 3,3 \quad \text{и} \quad \bar{x}_2 = 2,48 \quad .$$

Известно, что генеральные дисперсии примерно равны, т. е. $D_{ГЕН_1} = D_{ГЕН_2}$.

При уровне значимости $\alpha \leq 0,05$ проверить нулевую гипотезу $H_0: M(X_1) = M(X_2)$ если

$$t_{НАБЛ} = 3,27 \quad .$$

Решение.

$$t_{\text{КРИТ}}(\alpha \leq 0,05, f = n_1 + n_2 - 2 = 5 + 6 - 2 = 9) = 2,26.$$

$$t_{\text{НАБЛ}} > t_{\text{КРИТ}}(\alpha, f) \Rightarrow \text{отвергаем } H_0$$

Вывод: выборочные средние различаются
значимо $\alpha \leq 0,05$

4. Проверка гипотез о законах распределения

Во многих практических задачах закон распределения случайных величин заранее не известен, и надо выбрать модель, согласующуюся с результатами наблюдений.

Выдвигают нулевую гипотезу: неизвестная функция распределения исследуемой случайной величины X распределена по некоторому теоретическому закону, например, по нормальному закону

$$H_0 : F(x) = F_{TEOP}(x)$$

В качестве этой теоретической модели может быть рассмотрен любой закон, например, экспоненциальный или биномиальное распределение.

Это определяется сущностью изучаемого явления, а также результатами предварительной обработки наблюдений: формой графика распределения, соотношениями между выборочными данными.

Выдвигается альтернативная гипотеза, что данная генеральная совокупность не распределена по закону $F_{ТЕОР}(x)$:

$$H_1 : F(x) \neq F_{ТЕОР}(x)$$

Задается уровень значимости, например,

$$\alpha \leq 0,05$$

Если хотим проверить, согласуются эмпирические данные с нашим гипотетическим предположением относительно теоретической функции распределения или нет, то используем критерий согласия.

Критерий согласия – это критерий проверки гипотезы о предполагаемом законе неизвестного распределения.

Рассмотрим один из них, использующий распределение χ^2 и получивший название **критерий согласия Пирсона**.

Применим критерий χ^2 к проверке нулевой гипотезы, что генеральная совокупность распределена нормально.

Критерий предполагает, что результаты наблюдений сгруппированы в вариационный ряд и разбиты на классы.

По выборке объема n построим эмпирическое распределение $F_{ЭМП}(x)$:

варианты: x_1, x_2, \dots, x_k ;

эмпирические частоты: n_1, n_2, \dots, n_k ;

и сравним его с предполагаемым теоретическим распределением, вычисленным в предположении нормального закона распределения.

Теоретические частоты: n'_1, n'_2, \dots, n'_k .

То есть фактически $H_0 : n_{ЭМП} = n'_{ТЕОР}$

В качестве критерия проверки нулевой гипотезы примем случайную величину:

$$\chi^2_{НАБЛ} = \sum_{i=1}^k \frac{(n_{ЭМП} - n'_{ТЕОР})^2}{n'_{ТЕОР}},$$

где k – число классов.

Из таблиц находим $\chi^2_{КРИТ}(\alpha \leq 0,05; f = k - 3)$

Сравниваем, если $\chi^2_{НАБЛ} < \chi^2_{КРИТ}(\alpha, f) \Rightarrow H_0$

- расхождение теоретических и эмпирических частот незначимое. Следовательно, данные наблюдений согласуются с гипотезой о нормальном законе распределения генеральной совокупности.

Пример.

При уровне значимости $\alpha \leq 0,05$ проверить гипотезу о нормальном распределении генеральной совокупности, если известны эмпирические и теоретические частоты.

эмпирические частоты:

6 13 38 74 106 85 30 14;

теоретические частоты:

3 14 42 82 99 76 37 13.

Решение.

$$\chi^2_{\text{НАБЛ}} = 7,19$$

Найдем $\chi^2_{\text{КРИТ}}(\alpha \leq 0,05, f = 8 - 3 = 5) = 11,1$

Сравниваем: $\chi^2_{\text{НАБЛ}} < \chi^2_{\text{КРИТ}}(\alpha, f) \Rightarrow H_0$
- расхождение теоретических и эмпирических частот незначимое.

Следовательно, данные наблюдений согласуются с гипотезой о нормальном законе распределения генеральной совокупности.