## Лекция 11. ЭЛЕМЕНТЫ МАТЕМАТИЧЕСКОЙ СТАТИСТИКИ

Основная цель теории вероятностей - исследование имеющейся вероятностной модели. Математическая статистика ставит перед собой обратную задачу — построение вероятностно-статистической модели по результатам наблюдений.

<u>Из истории математической статистики.</u> Как наука математическая статистика возникла в XXв., но отдельные задачи рассматривались еще в XVIIв. Математическая статистика развивалась параллельно с теорией вероятностей. Дальнейшее развитие математической статистики (вторая половина XIX — начало XX вв.) связано с именами Чебышева, Маркова, Ляпунова, а также Гаусса, Кетле, Гальтона, К. Пирсона и др. В XXв. большой вклад в развитие математической статистики внесли советские математики Романовский, Слуцкий, Колмогоров, Смирнов, а также английские Стьюдент, Р. Фишер, Э. Пирсон и американские Нейман, Вальд.

Основные задачи математической статистики. Основная задача математической статистики — получение выводов о массовых явлениях и процессах по данным наблюдений над ними или экспериментов. Эти выводы относятся не к отдельным испытаниям, а являются утверждениями об общих характеристиках изучаемого явления — вероятностях, законах распределения и их параметрах и др. — в предположении постоянства условий, порождающих явление.

Содержание математической статистики составляет разработка приемов статистического наблюдения и анализа статистических данных. Исходный материал для статистического исследования реального явления - набор результатов наблюдений над ним или результатов специальных испытаний. Укажем некоторые основные вопросы, которые при этом возникают.

- 1. Оценка неизвестной вероятности случайного события.
- 2. Определение неизвестной функции распределения. В результате п независимых испытаний над случайной величиной  $\xi$  получены ее значения  $x_1, x_2, ..., x_n$ . Требуется определить, хотя бы приближенно, неизвестную функцию распределения  $F_{\xi}(x)$  величины  $\xi$ .
- 3. О пределение неизвестных параметров распределения. Случайная величина  $\xi$  имеет функцию распределения определенного вида, зависящую от k параметров, значения которых неизвестны. На основании последовательных наблюдений величины  $\xi$  нужно найти значения этих параметров.
- 4. П р о в е р к а с т а т и с т и ч е с к и х г и п о т е з. На основании некоторых соображений можно считать, что функция распределения случайной величины  $\xi$  есть  $F_{\xi}(x)$ ; спрашивается, совместимы ли наблюденные значения с гипотезой, что  $\xi$  действительно имеет распределение  $F_{\xi}(x)$ ?

Задачами 1-4 не исчерпываются основные проблемы математической статистики. Современную математическую статистику определяют как науку о принятии решений в условиях неопределенности.

<u>Генеральная совокупность и выборка.</u> В практике статистических наблюдений различают два основных вида наблюдений: сплошное, когда изучаются все элементы интересующей исследователя совокупности, и выборочное, когда изучается часть элементов. Вся подлежащая изучению совокупность называется *генеральной совокупностью*. Та часть объектов, которая отобрана для непосредственного изучения из генеральной совокупности, называется *выборкой*.

В математической статистике генеральная совокупность - это совокупность всех мыслимых наблюдений, которые могли бы быть произведены при данном реальном

комплексе условий и аналогично понятию случайной величины, так как полностью обусловлено комплексом условий.

О пределение всемножество значений подлежащей исследованию случайной величины называется *генеральной совокупностью*.

О п р е д е л е н и е. Последовательность наблюдений  $x_1, x_2, ..., x_n$  называется (случайной) выборкой объема n, если  $x_1, x_2, ..., x_n$  получены как независимые реализации некоторой случайной величины  $\xi$  с функцией распределения  $F_{\xi}(x)$ .

При этом также говорят, что  $x_1, x_2,..., x_n$  есть выборка из генеральной совокупности  $\xi$ . С теоретико-вероятностной точки зрения случайная выборка  $x_1, x_2,...,x_n$  может рассматриваться как последовательность независимых случайных величин, имеющих одну и ту же функцию распределения  $F_{\xi}(x)$ .

О пределение. Последовательность выборочных значений, записанных в порядке возрастания, называется *вариационным рядом*.

Пусть из генеральной совокупности извлечена выборка объема п:

$$\underbrace{Z_{1},Z_{1},...,Z_{1}}_{n_{1}\text{ pa3}};\underbrace{Z_{2},Z_{2},...,Z_{2}}_{n_{2}\text{ pa3}};...;\underbrace{Z_{k},Z_{k},...,Z_{k}}_{n_{k}\text{ pa3}}(n_{1}+n_{2}+...+n_{k}=n).$$

О п р е д е л е н и е. Различные выборочные значения  $z_1, z_2,...,z_k$  называются вариантами, числа  $n_1, n_2,...n_k$  – их частот равна объему выборки), а числа  $w_1 = n_1/n$ ,  $w_2 = n_2/n,...$ ,  $w_k = n_k/n$  – их относительными частот равна единице).

Zi	$\mathbf{z}_1$	$\mathbf{Z}_2$		Zk
$n_{\rm i}$	$n_1$	$n_2$	•••	$n_k$

называется статистическим рядом.

Иногда статистическим рядом (или рядом относительных частот) называют таблицу вида

Zi	$\mathbf{z}_1$	<b>Z</b> 2	 Zk
Wi	$\mathbf{w}_1$	W2	 Wk

3 а м е ч а н и е. Последний статистический ряд в математической статистике – аналог ряда распределения дискретной случайной величины генеральной совокупности), где вместо вероятностей  $p_i$  значений случайной величины стоят относительные частоты  $w_i$  вариант.

О п р е д е л е н и е. Функцию  $F_n(x) = m/n$ , где n — объем выборки, а m — число значений  $x_i$  в выборке, не превосходящих x, называют *эмпирической функцией распределения*.

В отличие от  $F_n(x)$ , полученной на основании выборочных значений генеральной совокупности, функцию распределения всей генеральной совокупности  $F_\xi(x)$  называют *теоретической функцией распределения*. Различие этих функций в том, что  $F_\xi(x)$  есть вероятность события  $\{\xi \leq x\}$ , а  $F_n(x)$  — относительная частота того же события. Из теоремы Бернулли следует, что при больших п числа  $F_n(x)$  и  $F_\xi(x)$  мало отличаются друг от друга в том смысле, что для любого положительного числа  $\varepsilon$  выполняется:  $P\{|F_n(x)-F_\xi(x)|<\varepsilon\}\to 1$  при  $\varepsilon$ . Это подтверждается и тем, что эмпирическая функция распределения обладает свойствами, аналогичными свойствам теоретической. Некоторые из них:

- 1)  $0 \le F_n(x) \le 1$ .
- 2)  $F_n(x)$  неубывающая функция, ее график имеет ступенчатый вид.
- 3) Если  $z_1$  наименьшая варианта, то  $F_n(x) = 0$  при  $x < z_1$ ; если  $z_k$  наибольшая варианта, то  $F_n(x) = 1$  при  $x \ge z_k$ .

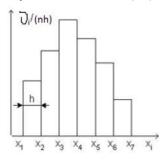
Одно из основных понятий математической статистики – оценка, или статистика.

О п р е д е л е н и е. Пусть  $x_1, x_2, ..., x_n$  – выборка. Любая функция от выборочных значений  $x_1, x_2, ..., x_n$  называется *оценкой* (той или иной характеристики генеральной совокупности).

Итак, эмпирическая функция распределения  $F_n(x)$  может служить оценкой по выборке для (неизвестной) теоретической функции распределения  $F_{\xi}(x)$  генеральной совокупности.

Если изучаемая случайная величина  $\xi$  (генеральная совокупность) непрерывна, то для оценки ее неизвестной плотности вероятности р $\xi$ (x) целесообразно использовать т.н. *гистограмму*, которая строится следующим образом.

- Упорядоченная по возрастанию выборка разбивается на интервалы (их число N можно определить, напр., по формуле Стерджесса N ≈ 1+3.322lg n, где n объем выборки, эта формула дает хорошие результаты, если n достаточно велико, распределение близко к нормальному, и используются равные интервалы).
- 2) Длина h каждого интервала (в случае их одинаковых длин) находится как  $h=(x_{max}-x_{min})/N$ , где  $x_{max}$  ( $x_{min}$ ) наибольшее (наименьшее) выборочное значение.
- 3) По оси х откладываются границы интервалов  $x_{min}$ ,  $x_{min}$ +h;  $x_{min}$ +2h; ...;  $x_{min}$ +Nh= $x_{max}$ , а над самими интервалами строятся прямоугольники так, чтобы площадь каждого из них представляла собой относительную частоту распределения:  $hl_i = n_i/n$ , где  $l_i$  высота i-го прямоугольника,  $n_i$  частота попадания выборочных значений в i-й интервал (если выборочное значение попадает на границу i-го и (i+1)-го интервалов, то считается, что оно принадлежит (i+1)-му интервалу. Имеем:  $l_i = n_i/(nh)$ .



Заметим (аналог нормировки плотности:  $\int\limits_{-\infty}^{+\infty}p_{\xi}(x)dx=1$ ), что общая площадь прямоугольников составляет 100% распределения:  $S=\sum_{i=1}^{N}hn_{i}/(nh)=1$ .

Выборочное среднее и выборочная дисперсия. Пусть имеется (повторная) выборка объема и из генеральной совокупности, т.е. N независимых одинаково распределенных случайных величин  $x_1, x_2, ..., x_n$  с функцией распределения  $F_\xi(x)$ .

О пределение. Выборочным средним называется число  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  (для статистического ряда  $\bar{x} = \frac{1}{n} \sum_{i=1}^k z_i n_i$  ).

Запишем последнюю формулу в виде  $\bar{x}=z_1n_1/n+z_2n_2/n+...+z_kn_k/n=z_1w_1+z_2w_2+...+z_kw_k$ . При больших n (в силу теоремы Бернулли) относительные частоты значений  $z_i$  близки к их вероятностям  $p_i$ :  $w_i\approx p_i$ , и мы получаем:  $\bar{x}\approx z_1p_1+z_2p_2+...+z_kp_k=M\xi$ , где  $\xi$  - (дискретная) случайная величина, представленная выборкой. Итак, в качестве оценки математического ожидания  $M\xi$  генеральной совокупности можно взять выборочное среднее  $\bar{x}$ .

Выборочное среднее, найденное по данным одной выборки, есть определенное число. Если извлекать другие выборки того же объема из генеральной совокупности, то выборочное среднее будет меняться от выборки к выборке. Таким образом, выборочное среднее можно рассматривать как случайную величину, а следовательно, говорить о ее распределении (теоретическом и эмпирическом) и их числовых характеристиках.

В качестве оценки для дисперсии генеральной совокупности по выборочным данным берут т. н. выборочную дисперсию.

О п р е д е л е н и е. Выборочной дисперсией называется число  $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \overline{x})^2$ . Если выборка записана в виде статистического ряда, то выборочная дисперсия вычисляется как  $s^2 = \frac{1}{n} \sum_{i=1}^k (z_i - \overline{x})^2 n_i$ 

Эта формула – аналог теоретико-вероятностной формулы для вычисления дисперсии (дискретной) случайной величины.

Часто бывает легче вычислять выборочную дисперсию по формуле  $s^2=\overline{x^2}$  -  $\overline{x}^2$ . Действительно,  $s^2=\frac{1}{n}\sum_{i=1}^k(z_i-\overline{x})^2n_i=\frac{1}{n}\sum_{i=1}^kz_i^2n_i$  -  $2\,\overline{x}\,\frac{1}{n}\sum_{i=1}^kz_in_i+\overline{x}^2\frac{1}{n}\sum_{i=1}^kn_i=\overline{x^2}$  -  $2\,\overline{x}^2+\overline{x}^2=\overline{x}^2$  -  $2\,\overline{x}^2$  -  $2\,$ 

3 а м е ч а н и е. Обычно (ниже будет показано, почему) в практике статистических вычислений вместо  $s^2$  используют оценку  $s_1^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \overline{x})^2$  или (для статистического ряда)  $s_1^2 = \frac{1}{n-1} \sum_{i=1}^k (z_i - \overline{x})^2 n_i$ , называемую исправленной выборочной дисперсией.

О п р е д е л е н и е. Выборочным среднеквадратическим отклонением (стандартным отклонением) называется число  $s=\sqrt{s^2}~(s_1=\sqrt{s_1^2}~).$ 

Выборочная ковариация и выборочный коэффициент корреляции. Пусть при проведении некоторого опыта наблюдаются две случайные величины  $\xi$  и  $\eta$ . Тогда n независимых повторений опыта дадут n пар наблюдавшихся значений  $(x_1,y_1)$ ,  $(x_2,y_2)$ , ..., $(x_n,y_n)$   $(x_i$  – выборочные значения  $\xi$ ,  $y_i$  –выборочные значения  $\eta$ ).

О п р е д е л е н и е. Выборочной ковариацией случайных величин  $\xi$  и  $\eta$  называется число

$$S_{\xi,\eta} = \tfrac{1}{n} \sum_{i=1}^n (x_i - \overline{x}) (y_i - \overline{y}) = \tfrac{1}{n} \sum_{i=1}^n x_i y_i - \tfrac{1}{n} \sum_{i=1}^n x_i \, \frac{1}{n} \sum_{i=1}^n y_i \; .$$

Выборочная ковариация служит оценкой по выборке ковариации случайных величин  $\xi$  и  $\eta$  (генеральной ковариации) и является выборочной характеристикой их линейной связи.

О п р е д е л е н и е. Выборочным коэффициентом корреляции случайных величин  $\xi$  и  $\eta$  называется число  $r_{\xi,\eta} = S_{\xi,\eta}/(s_{\xi}s_{\eta})$ , где  $s_{\xi}$ ,  $s_{\eta}$  - выборочные среднеквадратические отклонения случайных величин  $\xi$  и  $\eta$  (соответственно).

Выборочный коэффициент корреляции служит оценкой по выборке коэффициента корреляции случайных величин  $\xi$  и  $\eta$  (генерального коэффициента корреляции) и является выборочной характеристикой направления и тесноты стохастически-линейной линейной связи между случайными величинами  $\xi$  и  $\eta$ .