

## Математическая статистика

## Вопрос 1. Предмет математической статистики и ее основные задачи

В теории вероятностей, если мы изучаем случайную величину X, ее закон распределения считается заданным, и мы можем достоверно ответить на любой вопрос, касающийся данной случайной величины В математической

статистике ситуация прямо противоположная — мы ничего не знаем о законе распределения изучаемой случайной величины X. У нас имеются только некоторые ее наблюдения или измерения. Понятно, что по конечному числу наблюдений невозможно достоверно сделать какие-либо выводы об изучаемой случайной величине. Ясно также, что чем больше таких наблюдений, тем более надежными будут наши приближенные выводы

В этом состоит основная особенность математической статистики— она не определяет достоверно закономерности поведения изучаемых случайной явлений, а оценивает их с той или иной степенью достоверности. Но при неограниченном увеличении числа наблюдений выводы математической статистики становятся практически достоверными.

Поэтому содержание этой дисциплины — как и сколько сделать наблюдений и как их обработать, чтобыответить на интересующий нас вопрос о случайном явлении с требуемой степенью достоверности. Итак, установление закономерностей, которым подчинены массовые случайные явления основано на изучении статистических данных —

Митематическая статистика решаетдве главные задачи:

результатах наблюдений.

- 1) указать способы сбора и группировки (если данных очень много) статистических сведений (результатов наблюдений);
- 2) разработать методы анализа собранных статистических данных в зависимости от целей исследования.

Вопрос 2. Выборка. Статистический ряд. Эмпирический закон распределения. Полигон и гистограмма



0000 × 000 + 00

Пусть требуется изучить совокупность однородных объектов относительно качественного или количественного признака, характеризующего эти объекты

Можно производить этот контроль <u>сплошным обследованием</u>, то есть измерять каждый из объектов совокупности. Но на практике сплошное обследование применяется редко:

- а) из-за очень большого числа объектов;
- б) из-за того, что иногда обследование заключается в физическом уничтожении, например, проверяем взрываемость гранат или проверяем на

крепость произведенную посуду и т.д.

В таких случаях производится случайный отбор ограниченного (небольшого) числа объектов, которые и подвергают изучению.

**Выборочной совокупностью** (выборкой) называется совокупность случайно отобранных однородных объектов.

*Генеральной совокупностью (ГС)* называется совокупность всех однородных объектов, из которых производится выборка.

**Объемом** совокупности (выборочной или генеральной) называется число объектов этой совокупности.

При наборе выборки можно поступать двояко: после того, как объект отобран и над ним произведено наблюдение, он может быть возвращен либо не возвращен в генеральную совокупность. В связи с этим выборки подразделяются на *повторные* и *бесповторные*.

Для того, чтобы по данным выборки можно было достаточно уверенно судить об интересующем нас признаке ГС, необходимо, чтобы объекты выборки правильно его представляли. Это требование коротко формулируется так: выборка должна быть *репрезентативной* (представительной).

Способы отбора выборки:

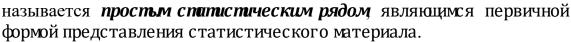
- 1. Отбор, не требующий расчленения ГС на части:
- а) простой случайный бесповторный;
- б) простой случайный повторный.
- 2. Отбор, при котором ГС разбивается на части (если объем ГС слишком большой):
- а) типический отбор. Объекты отбираются не из всей ГС, а из ее «типичных» частей. Например, цех из тридцати станков производит одну и ту же деталь. Тогда отбор делается по одной или по две детали с каждого станка в случайные моменты времени;
- б) механический отбор. Например, если нужно выбрать 5% деталей, то выбирают не случайно, а каждую двадцатую деталь;
- в) серийный отбор. Объекты выбирают не по одному, а сериями.

Итак, пусть из ГС значений некоторого количественного признака произведена вьборка объема N:  $X = \{x_1, \dots, x_N\}$ .



Таблица вида

№	1	2	•••	N
X	$\mathbf{x}_1$	$\mathbf{x}_2$	•••	$X_N$



Если данные таблицы, называемые *вариантами*, упорядочены по возрастанию, то выборка  $X = \{x_1, \dots, x_N\}$  называется *вариационным рядом* 

**Размах** выборки — это длина основного интервала [ $x_{min}$ ;  $x_{max}$ ], в который попадают все значения выборки. Вынисляется размах выборки следующим образом:  $d = x_{max} - x_{min}$ . Затем по формуле  $k = 1 + 3,322 \cdot lgN$ , определяется число k. Данное число задает количество подынтервалов (классов), на которые разбиваем основной интервал. Длиныһ подынтервалов и их границыа $_i$  ( $i = \overline{0,k}$ ) вынисляются по формулам:  $h = \frac{d}{k}$ ,  $a_0 = x_{min}$ ,  $a_1 = a_0 + h$ , ...,  $a_i = a_{i-1} + h$ 

Далее находятся **часпоты**  $m_j$   $(j=\overline{1,k})$  и **относительнье часпоты** $\mu_j=\frac{m_j}{N}$   $(j=\overline{1,k})$  попадания значений вьборки X в j-й подынтервал. Причем для частот должно выполняться равенство  $\sum_{j=1}^k m_j = N$ , а для относительных

частот соответственно  $\sum_{j=1}^{k} \mu_{j} = 1$ .

Огносительные частоты иногда называют *эмпирическими вероятностями*.

Результаты проведенных расчетов сводятся в таблицу, которая называется **интервальнымвариационнымрядом** 

X	$[a_0;a_1)$	$[a_1;a_2)$	•••	$[a_{k-1};a_k]$
m	$m_1$	$m_2$		$m_k$
μ	$\mu_1$	$\mu_2$		$\mu_k$

Далее находятся середины подыттервалов:  $\bar{x_i} = \frac{a_i + a_{i-1}}{2}$  и составляется таблица, которая называется дискрепный вариационный ряд:

X	$\overline{x_1}$	$\overline{X_2}$	•••	$\overline{\mathbf{x}_{k}}$
m	$m_1$	$m_2$	• • •	$m_k$
μ	$\mu_1$	$\mu_2$	• • •	$\mu_k$

Вариационный ряд распределения является оценкой теоретического ряда распределения и сходится к нему по вероятности. Поскольку ряд распределения является одной из формзадания закона распределения дискретной случайной величины, то мы получили эмпирический закон распределения исследуемой дискретной случайной величины

Сгруппированные данные вариационного ряда несут в себе меньше информации, чем выборочные, так как в них теряется информация о



0000 X 000 + +

порядке следования вьборочньк значений. При группировке также фактически происходит округление наблюдаемых значений выборки внутри j-го класса

(подыттервала) до значения  $x_j$ , что приводит к потере информации о распределении исследуемой случайной величины внутри каждого класса. Это распределение в дальнейшем предполагается равномерным Преимуществом же сгруппированных данных является их компактность и большая наглядность.

В целях визуального изучения полученных вариационных рядов пользуются различными способами их графического изображения. К ним относятся гистограмма и полигон.

**Гистограммой** относительных частот называется ступенчатая фигура, состоящая из прямоугольников, основаниями которых служат частичные

подытервалы длины h , а высоты равны числам  $\frac{\mu_j}{h}$  (плотности вероятностей)  $(j=\overline{1,k})$ .

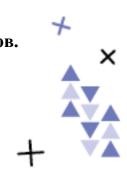
Площадь гистограммы относительных частот равна единице. Аналогичным образом строится гистограмма частот. Площадь гистограммычастот равна объему выборки.

**Полигоном** относительных частот называется ломаная, отрезки которой соединяют точки  $(\overline{x_1}; \mu_1), ..., (\overline{x_k}; \mu_k)$ . Полигон относительных частот есть визуальное представление эмпирического закона распределения выборки.

Любая функция выборки  $\varphi = \varphi(x_1, ..., x_N)$  называется *статистикой*. Статистика является случайной величиной, так как на различных реализациях выборки она получает различные наблюдаемые значения. Статистиками являются: частоты  $m_j$ , границы классов  $a_j$  и их середины  $\overline{x_j}$ , размах. Вариационный ряд распределения также является статистикой. Из определения статистики следует, что любая функция от статистик также является статистикой, поэтому статистикой является любая функция от сгруппированных данных.

Статистики служат для оценки любых характеристик изучаемой случайной величины: вероятностей случайных событий, связанных с изучаемой величиной, ее числовых характеристик, параметров закона распределения и так далее. Изучение статистик на основе теории вероятностей есть теоретическое ядро математической статистики.

Вопрос 3. Статистические оценки генеральных параметров. Точечные и интервальные оценки



0 0 0 0 X 0 0 0 7

×

Пусть  $\Theta_{\Gamma}$  — некоторый параметр  $\Gamma$ С, который невозможно вычислить. Но знать его значение (хотя бы приближенное, оценочное) надо! Поэтому по выборочным данным производят расчет статистических оценок данного генерального параметра.

Оценки параметров подразделяются на точечные и интервальные.

**Точечной** называется статистическая оценка генерального параметра  $\Theta_{\Gamma}$ , которая определяется одним числом  $\widetilde{\Theta}_{B}$ .

**Интервальной** называется оценка генерального параметра  $\Theta_{\Gamma}$ , которая определяется двумя числами  $\widetilde{\Theta}_{\rm B}$  и  $\widetilde{\widetilde{\Theta}}_{\rm B}$  — концами интервала, покрывающего оцениваемый генеральный параметр  $\Theta_{\Gamma}$ .

Для того, чтобы точечная оценка давала «хорошие» приближения оцениваемого параметра, она должна быть: несмещенной, эффективной, состоятельной.

**Несмещенной** называют такую точечную оценку  $\widetilde{\Theta}_{\rm B}$ , математическое ожидание которой равно оцениваемому генеральному параметру при любом

объеме выборки:  $M(\theta_B) = \theta_\Gamma$ , в противном случае оценка называется *смещенной*.

**Эффективной** называется точечная оценка  $\widetilde{\Theta}_{\rm B}$ , которая (при заданном объеме выборки) имеет наименьшую возможную дисперсию:  $M[(\theta_{\rm B} - \theta_{\Gamma})^2] \to min$ .

**Состоятельной** называется точечная оценка  $\widetilde{\Theta}_{\rm B}$ , которая (сувеличением объема выборки) стремится по вероятности к оцениваемому параметру  $\Theta_{\Gamma}$ :  $\forall \delta > 0$   $\lim_{N \to \infty} P(|\theta_{\Gamma} - \theta_{\rm B}| < \delta) = 1$ .

Несмеценной оценкой генеральной средней (генерального математического ожидания  $M_{\Gamma}[x]$ ) служит выборочная средняя (выборочное математическое ожидание):  $M_{\rm B}[x] = \sum_{j=1}^k \overline{x_j} \cdot \mu_j$ , где  $\overline{x_j}$  и  $\mu_j$  — данные дискретного вариационного ряда.

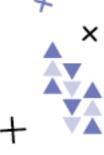
Кроме того,  $M_{\rm B}[x]$  является состоятельной оценкой. Если случайная величина X подчинена нормальному закону распределения, то  $M_{\rm B}[x]$  является и эффективной оценкой.

Смещенной оценкой генеральной дисперсии  $D_{\Gamma}[x]$  служит выборочная дисперсия:  $D_{\mathbb{B}}[x] = \sum_{j=1}^k \left(\overline{x_j} - M_{\mathbb{B}}[x]\right)^2 \cdot \mu_j$ .

Поскольку  $D_{\scriptscriptstyle \rm B}[x]$  является смещенной оценкой, то ее «исправляют» следующим образом:  $\sigma_{\scriptscriptstyle \rm B}^2=\frac{N}{N-1}\cdot D_{\scriptscriptstyle \rm B}[x].$ 

Полученная оценка  $\sigma_{\scriptscriptstyle B}^2$  — это состоятельная несмещенная выборочная дисперсия, а  $\sigma_{\scriptscriptstyle B}$  — выборочное среднее квадратичное отклонение.

При выборке малого объема (N<30) точечная оценка может значительно отличаться от оцениваемого генерального параметра, то есть приводить к грубым ошибкам. Поэтому при небольшом объеме выборки следует пользоваться интервальными оценками.



0000 X 000 ⊁ Пусть найленная

Пусть найденная (по данным выборки) статистическая оценка  $\Theta_{\rm B}$  является оценкой неизвестного генерального параметра  $\Theta_{\Gamma}$ . Ясно, что  $\Theta_{\rm B}$  тем точнее определяет  $\Theta_{\Gamma}$ , чем меньше значение разности  $|\Theta_{\Gamma} - \Theta_{\rm B}|$ . То есть при  $|\Theta_{\Gamma} - \Theta_{\rm B}| < \delta \ (\delta > 0)$  чем меньше  $\delta$ , тем оценка  $\Theta_{\Gamma}$  точнее. Значит, положительное число  $\delta$  характеризует *точность* оценки.

**Надежностью** (доверительной вероятностью) оценки  $\Theta_{\rm B}$  называется вероятность  $\gamma$ , с которой осуществляется событие  $|\Theta_{\Gamma} - \Theta_{\rm B}| < \delta$ , то есть  $\gamma = P(|\Theta_{\Gamma} - \Theta_{\rm B}| < \delta)$ .

Обынно надежность оценки (доверительная вероятность  $\gamma$ ) задается. Причем в качестве  $\gamma$  берут число, близкое к единице (0,95; 0,99; 0,999). **Доверительным** называется интервал, который с заданной надежностью  $\gamma$  покрывает оцениваемый генеральный параметр.

Vв формулы, определяющей  $\gamma$ , если раскрыть модуль, получается  $P(-\delta < \Theta_{\Gamma} - \Theta_{B} < \delta) = \gamma$  или  $P(\Theta_{B} - \delta < \Theta_{\Gamma} < \Theta_{B} + \delta) = \gamma$ .

Тогда интервал  $(\Theta_B - \delta; \Theta_B + \delta)$  и есть доверительный интервал. Из общих соображений ясно, что длина доверительного интервала будет зависеть от объема выборки N и доверительной вероятности  $\gamma$ .

Для оценки математического ожидания  $M_{\Gamma}[x]$  нормально распределенной ГС X по выборочной средней  $M_{\rm B}[x]$  при известном среднем квадратическом отклонении  $\sigma_{\Gamma} = \sqrt{D_{\Gamma}[x]}$  служит доверительный интервал:  $M_{\rm B}[x] - t \cdot \frac{\sigma_{\Gamma}}{\sqrt{N}} < M_{\Gamma}[x] < M_{\rm B}[x] + t \cdot \frac{\sigma_{\Gamma}}{\sqrt{N}}$ , где  $t \cdot \frac{\sigma_{\Gamma}}{\sqrt{N}} = \delta$  — точность оценки; N — объем выборки; t — такое значение функции Лапласа  $\Phi(t)$ , при котором  $\Phi(t) = \frac{\gamma}{2}$ .

Для оценки математического ожидания  $M_{\Gamma}[x]$  нормально распределенной ГС X по выборочной средней  $M_{\rm B}[x]$  при неизвестном среднем квадратическом отклонении  $\sigma_{\Gamma} = \sqrt{D_{\Gamma}[x]}$  (при объеме выборки  $N{>}30$ ) служит доверительный интервал:  $M_{\rm B}[x] - t_{\gamma} \cdot \frac{\sigma_{\rm B}}{\sqrt{N}} < M_{\Gamma}[x] < M_{\rm B}[x] + t_{\gamma} \cdot \frac{\sigma_{\rm B}}{\sqrt{N}}$ , где  $t_{\gamma}$  такое значение функции Лапласа  $\Phi(t)$ , при котором  $\Phi(t) = \gamma$ .

Для оценки среднего квадратического отклонения  $\sigma_{\Gamma}$  нормально распределенной ГС X с доверительной вероятностью  $\gamma$  служат доверительные интервалы  $\sigma_{\rm B}(1-q) < \sigma_{\Gamma} < \sigma_{\rm B}(1+q) \ \, ({\rm при} \ q < 1) \ \, 0 < \sigma_{\Gamma} < \sigma_{\rm B}(1+q) \ \, ({\rm при} \ q > 1) \ \,$  где  ${\rm q}$  – находится из таблицы при заданном  ${\rm N}$  и  ${\rm \gamma}$ .

3амечание 1. Оценку  $|\Theta_{\Gamma} - \Theta_{\rm B}| < t \cdot \frac{\sigma_{\Gamma}}{\sqrt{N}}$  называют классической. Из формулы  $\delta = t \cdot \frac{\sigma_{\Gamma}}{\sqrt{N}}$ , определяющей точность классической оценки, можно сделать выводы

