

1. Данные в экономике. Классификация данных.

2. Генеральная совокупность и выборка. Суть выборочного метода.

Генеральная совокупность Ω – это совокупность всех подлежащих изучению объектов или явлений. В некоторых задачах генеральную совокупность рассматривают как случайную величину X .

Выборочная совокупность (или выборка) Ω^{\sim} - совокупность случайно отобранных объектов из генеральной совокупности.

- Число объектов совокупности называется объемом совокупности
- Объем генеральной совокупности – N
- Объем выборки – n

Основной метод математической статистики – выборочный. Его суть:

- Выборочный метод – метод матстатистики, где на основе изучения выборки делается заключение обо всей генеральной совокупности.
- Теоретической основой применения выборочного метода является Закон Больших Чисел: при неограниченном увеличении объема выборки её характеристики сколь угодно близко приближаются к характеристикам генеральной совокупности.

3. Способы осуществления выборки. Условия репрезентативности выборки.

Чтобы выборка правильно представляла изучаемый признак генеральной совокупности, хорошо отражала пропорции генеральной совокупности, она должна быть репрезентативной (представительной). Выборка будет репрезентативной, если:

- Её осуществить случайно
- Все объекты генеральной совокупности имеют равные вероятности быть отобранными

Способы формирования выборки:

- Повторный (возвратный): объект после исследования возвращается в генеральную совокупность
- Бесповторный (безвозвратный): объект после исследования не возвращается в генеральную совокупность

4. Понятие вариационного ряда. Дискретные и интервальные статистические ряды: понятие, способы задания.

Статистический ряд – это ранжированный перечень вариантов x_i и соответствующих им весов (частот или частостей)

Общий вид статистического ряда частот/частостей

Значение признака x_i	x_1	x_2	...	x_k	
Частота n_i	n_1	n_2	...	n_k	$n_1+n_2+\dots+n_k=n$
Относительная частота ω_i	ω_1	ω_2	...	ω_k	$\omega_1+\omega_2+\dots+\omega_k=1$

Где k – число различных вариантов в ряду

Статистические ряды бывают дискретными и интервальными:

- Стат ряд называют дискретным, если любые его варианты отличаются друг от друга на постоянную величину. В таких рядах задаются точечные значения признака.

- Статистический ряд называется интервальным, если любые его варианты отличаются друг от друга на сколь угодно малую величину. Значения признака в них задаются в виде интервалов.

Если число значений признака X велико, то варианты разбивают на отдельные интервалы, т.е. проводят их группировку.

На практике обычно считают, что правильно составленный ряд распределения содержит от 5 до 15 частичных интервалов.

Рекомендуемое число интервалов вычисляется по формуле Стерджеса:

$$m = 1 + 3,322 \cdot \lg(n)$$

Ширина (величина) интервала h равна:

$$h = \frac{x_{max} - x_{min}}{m}$$

5. Эмпирическая функция распределения, её график и свойства.

Эмпирической (статистической) функцией распределения $F_n(x)$ называется функция, равная относительной частоте того, что признак (СВ X) примет значение меньше заданного действительного числа x , т.е. функция, определяющая, для каждого значения x частоту события $\{X < x\}$:

$$F_n(x) = \omega(X < x) = \frac{n_i}{n} = \frac{\text{число вариантов } x_i \text{ меньших } x}{\text{объем выборки}}$$

x - любое действительное число

Свойства эмпирической функции распределения:

- Значение $F_n(x)$ принадлежат отрезку $[0; 1]$
- $F_n(x)$ является неубывающей функцией
- $F_n(x) = 0$ при $x \leq x_{min}$
- $F_n(x) = 1$ при $x > x_{max}$

Для дискретного статистического ряда:

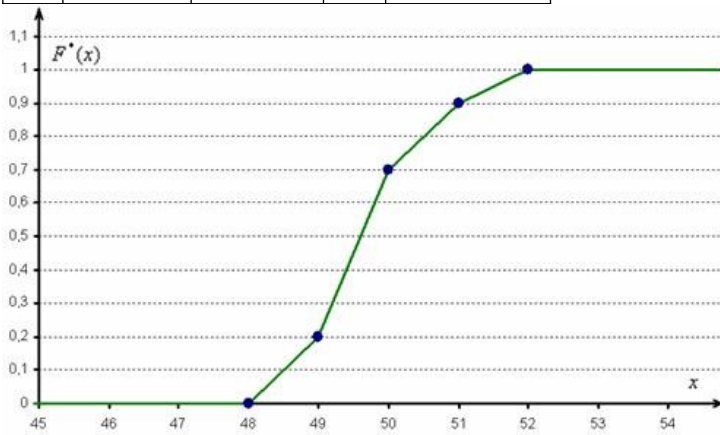
x_i	x_1	x_2	...	x_k
ω_i	ω_1	ω_2	...	ω_k



$$F_n(x) \begin{cases} 0 & \text{при } x \leq x_1 \\ \omega_1 & \text{при } x_1 < x \leq x_2 \\ \omega_1 + \omega_2 & \text{при } x_2 < x \leq x_3 \\ \omega_1 + \omega_2 + \omega_3 & \text{при } x_3 < x \leq x_4 \\ \dots & \\ 1 & \text{при } x > x_k \end{cases}$$

Для интервального статистического ряда:

x_i	$(x_1; x_2]$	$(x_2; x_3]$...	$(x_{k-1}; x_k]$
ω	ω_1	ω_2	...	ω_k
i				



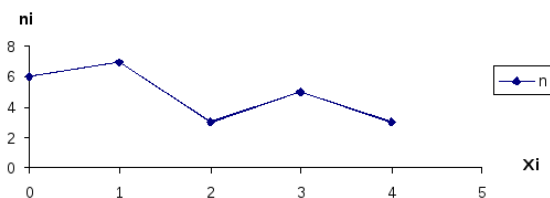
- $F(x)=0$ при $x \leq x_1$
- $F(x_2)=\omega_1$
- $F(x_3)=\omega_1+\omega_2$
- $F(x_4)=\omega_1+\omega_2+\omega_3$
- ...
- $F(x_k)=\omega_1+\omega_2+\dots+\omega_{k-1}+\omega_k=1$
- $F(x)=1$ при $x > x_k$

6. Графическое представление статистических рядов: полигон и гистограмма.

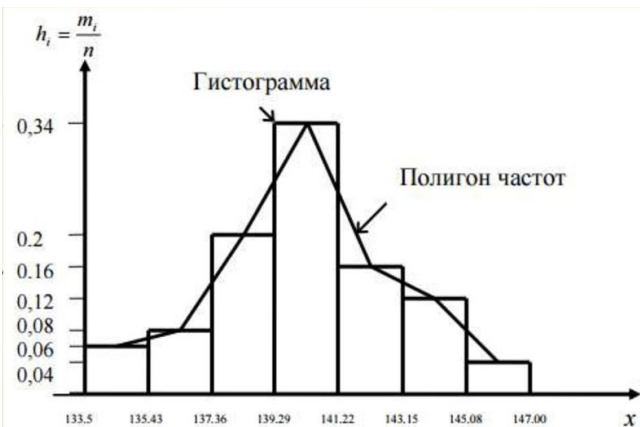
Полигон как правило служит для изображения дискретного статистического ряда. Полигон частот (или частостей) – это ломаная, отрезки которой соединяют точки с координатами $(x_i; n_i)$ или $(x_i; \omega_i)$, $i=1, 2, \dots, k$.

Варианты (x_i) откладывают на оси абсцисс, а частоты или частости – на оси ординат.

Полигон частот



Гистограмма (т.е. столбчатая диаграмма) служит только для изображения интервальных статистических рядов. Гистограмма частот или частостей – это ступенчатая фигура, состоящая из прямоугольников, основаниями которых служат интервалы $(x_i; x_{i+1})$, $i=1, 2, \dots, m$, а высоты равны частотам (n_i) или частостям (ω_i) .



7. Выбросы. Диаграмма размаха (ящик с усами).

8. Числовые характеристики выборочных распределений: выборочная средняя, мода, медиана, показатели вариации, показатели формы.

Числовые характеристики признака X , рассчитанные по выборке, называются выборочными характеристиками этого признака. Выборочные характеристики являются случайными величинами, а не константами.

Пусть дано стат распределение выборки объема n

А) Дискретный статистический ряд частот:

Значение признака x_i	x_1	x_2	...	x_k
Частота n_i	n_1	n_2	...	n_k

Б) Интервальный статистический ряд частот

$(x_{i-1}; x_i]$	$(x_1; x_2]$	$(x_2; x_3]$...	$(x_{k-1}; x_k]$
n_i	n_1	n_2	...	n_k

I Характеристика центра распределения (средние)

- Выборочная средняя – это среднее арифметическое всех значений выборки:
 - Простая – используется, когда данные наблюдения не сведены в вариационный ряд, либо когда все частоты равны 1 или одинаковы

$$\bar{x}_{выб} = \frac{\sum_{i=1}^n x_i}{n}$$

- Взвешенная – используется, когда частоты отличны друг от друга:

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

- Мода M_o вариационного (статистического) ряда – это вариант, которому соответствует наибольшая частота.
 - Для дискретного вариационного ряда мода равна значению варианты с наибольшей частотой
 - Мода интервального ряда определяется по формуле:

$$M_o = x_{M_o} + \frac{h * n_{M_o} - n_{M_o-1}}{(n_{M_o} - n_{M_o-1}) + (n_{M_o} - n_{M_o+1})}$$

Где x_{M_o} - нижняя граница модального интервала

h - ширина интервала

n_{M_o} - частота модального интервала

Модальным считается интервал, которому соответствует наибольшая частота

- Медианой M_e вариационного (статистического) ряда называется значение признака, приходящееся на середину ранжированного ряда наблюдений.

- Для дискретного вариационного ряда с нечетным числом членов медиана равна срединному варианту, а для ряда с четным числом членов, полусумме двух срединных вариантов, т.е.

Если $n=2k+1$, то медиана $Me=x_{k+1}$

Если $n=2k$, то медиана $Me=(x_k+x_{k+1})/2$

- Для интервального ряда медиана определяется по формуле:

$$Me = x_{Me} + \frac{h * \frac{n}{2} - S_{Me-1}}{n_{Me}}$$

x_{Me} -нижняя граница медианного интервала

h -ширина интервала

n_{Me} - частота медианного интервала

n -объем выборки

S_{Me-1} - сумма частот (накопленная частота) до медианного интервала

Медианным считается интервал, которому принадлежит значение признака с номером $n/2$ (если n – четное) или $(n+1)/2$ (если n -нечетное)

II Показатели вариации признака

- Выборочная дисперсия $D_{выб}$ – это среднее арифметическое квадратов отклонений значений признака от выборочной средней:

- Простая: $D_{выб}(X) = \frac{\sum_{i=1}^n (x_i - \bar{x}_{выб})^2}{n} = \overline{x^2} - \bar{x}_{выб}^2 = \frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}_{выб}^2$

- Взвешенная: $D_{выб}(X) = \frac{\sum_{i=1}^k (x_i - \bar{x}_{выб})^2 * i}{\sum_{i=1}^k i} = \overline{x^2} - \bar{x}_{выб}^2 = \frac{\sum_{i=1}^k x_i^2 * i}{\sum_{i=1}^k i} - \bar{x}_{выб}^2$

- Выборочное среднее квадратическое отклонение (стандартное отклонение) – это арифметическое значение корня квадратного из выборочной дисперсии – рассчитывается по формуле: $\sigma_{выб} = \sqrt{D_{выб}}$

- Размах вариации R – это число, равное разности между наибольшим и наименьшим вариантами ряда, т.е. $R = x_{max} - x_{min}$

- Выборочный коэффициент вариации равен процентному отношению выборочного СКО к выборочной средней, т.е. $V(X) = \frac{\sigma_{выб}}{\bar{x}_{выб}} * 100\%$

- Выборочный начальный момент порядка k :

$$\hat{a}_k = \frac{1}{n} \sum x_i^k \text{ или } \hat{a}_k = \frac{1}{n} \sum x_i^k * n_i$$

- Выборочный центральный момент порядка k :

$$\tilde{\mu}_k = \frac{1}{n} \sum (x_i - \bar{x}_{выб})^k \text{ или } \tilde{\mu}_k = \frac{1}{n} \sum (x_i - \bar{x}_{выб})^k * i$$

III показатели формы распределения признака:

- Квантили (ранговые характеристики): q -квантили; d - децили; p -перцентили или процентиля. – это значения признака, которые делят ранжированный ряд соответственно на 4, 10, 100 равных частей

- Асимметрия А. Показывает различия в вариации значений признака по одну и другую сторону от средней. $A_{выб} = \frac{\tilde{\mu}_3}{\tilde{\sigma}_{выб}^3}$
- Эксцесс Е. Это показатель островершинности или плосковершинности симметричного распределения по сравнению с нормальным распределением. $E_{выб} = \frac{\tilde{\mu}_4}{\tilde{\sigma}_{выб}^4} - 3$
- Асимметрия и эксцесс определяются для интервальных статистических рядов.

9. Понятие статистической оценки параметров распределения. Виды стат. оценок.

Статистическое оценивание – это определение приближенного значения неизвестного параметра генеральной совокупности по результатам наблюдения.

Параметр θ – это числовая характеристика генеральной совокупности.

Статистической оценкой ($\tilde{\theta}$) параметра тета называется его приближенное значение, зависящее от данных выбора.

Виды статистических оценок:

- Точечные
Точечной оценкой θ^* параметра θ называется числовое значение этого параметра, полученное по выборке объема n
- Интервальные
Интервальной оценкой параметра тета называется числовой интервал $(\theta_1; \theta_2)$, который с определенной вероятностью накрывает (содержит) неизвестное значение параметра.

10. Свойства статистических оценок: несмещенность, состоятельность, эффективность.

- Несмещенность (оценка совпадает с параметром θ). Стат оценка параметра θ называется несмещенной, если ее мат ожидание (среднее значение) равно оцениваемому параметру при любом объеме выборки $M(\tilde{\theta}) = \theta \forall n$. В противном случае мат ожидание не равно параметру θ оценка называется смещенной. В этом случае разность между мат ожиданием и параметром называют смещением или систематической ошибкой оценивания. Св-во несмещенности является желательным, но не обязательным.
- Эффективность – оценка обладает наименьшей степенью случайных отклонений от параметра θ .
Несмещенная оценка параметра θ называется эффективной, если она имеет наименьшую дисперсию среди всех возможных несмещенных оценок параметра, т.е. если ее дисперсия минимальна при заданном объеме выборки. Эффективность определяется отношением:

$$e = \frac{\sigma_{\tilde{\theta}_2}^2}{\sigma_{\tilde{\theta}}^2}$$

Где в числителе дисперсия эффективной оценки, а в знаменателе – дисперсия данной оценки.

Чем ближе e к 1, тем эффективней оценка.

Требование эффективности является желательным, но не обязательным.

- Состоятельность – оценка стремится к параметру θ с ростом объема выборки (n).
Оценка параметра называется состоятельной, если она удовлетворяет ЗБЧ (т.е. при $n \rightarrow \infty$ сходится по вероятности к оцениваемому параметру)
Другими словами, для любого положительного числа ε выполняется условие:

$$\forall \varepsilon > 0; \lim_{n \rightarrow \infty} P(|\tilde{\theta} - \theta| < \varepsilon) = 1$$

Свойства состоятельности является обязательным для любого правила оценивания.

11. Точечная оценка математического ожидания.

Пусть изучается СВ X с мат ожиданием a .

Несмещенной точечной оценкой мат ожидания $M(X) = \overline{x_{ген}}$ является выборочная средняя.

$$\overline{x_{ген}} \approx \overline{x_{выб}}$$

12. Точечные оценки дисперсии и среднего квадратического отклонения.

$D_{ген}(X) = \sigma_{ген}^2(X)$ – генеральная дисперсия имеет 2 точечные оценки:

- $D_{выб}(X) = \sigma_{выб}^2(X)$ – выборочная дисперсия, которая является смещенной оценкой
- S^2 – исправленная выборочная дисперсия (несмещенная)

$$S^2 = D_{выб} \frac{(X) * n}{n-1} \text{ (множитель Бесселя)}$$

$\sigma_{ген}$ имеет 2 точечные оценки:

- $\sigma_{выб}$ – выб СКО (стандартное отклонение) – смещенная $\sigma_{выб} = \sqrt{\sigma_{выб}^2}$
- S – исправленное выборочное СКО $S = \sqrt{S^2}$ (несмещенная)

13. Точечная оценка генеральной доли.

Пусть генеральная совокупность содержит N элементов, из которых M обладают некоторым признаком A . Тогда доля единиц, обладающих признаком A в генеральной доле обозначается

$p = \frac{M}{N}$, где N – объем генеральной совокупности, M – число элементов в генеральной совокупности с признаком A .

Пусть n – объем выборки; m – число элементов с признаком A в выборке. Тогда величина $\omega = \frac{m}{n}$ называется выборочной долей.

Генеральную долю можно рассматривать как вероятность появления события A в биномиальном законе распределения СВ X , а выборочную долю – как относительную частоту этого события в n независимых испытаниях.

Несмещенной и состоятельной точечной оценкой генеральной доли является выборочная доля ω ,

$$\text{т.е. } p \approx \omega = \frac{m}{n}$$

14. Методы нахождения точечных оценок: метод моментов, метод максимального правдоподобия.

Наиболее распространенные методы:

- Метод моментов (ММ)
- Метод максимального правдоподобия (ММП)
- Метод наименьших квадратов (МНК)

Метод моментов – это приравнивание теоретических моментов распределения (т.е. начальных α_k или центральных μ_k) СВ X к соответствующим эмпирическим (выборочным) моментам, найденным по выборке, с последующим решением полученного уравнения (системы уравнений).

Если теоретическое распределение задается одним параметром, то для нахождения точечной оценки неизвестного параметра необходимо решить одно уравнение вида:

$$M(X) = \overline{x_{выб}}$$

II Если теоретическое распределение задается двумя параметрами, то для их оценки необходимо решить систему из двух уравнений вида:

$$\begin{cases} M(X) = \overline{x_{\text{выб}}} \\ D(X) = \sigma_{\text{выб}}^2 \end{cases}$$

Достоинства и недостатки ММ:

- неприменим, когда моменты генеральной совокупности нужного порядка не существуют
- оценки, полученные данным методом, являются состоятельными, чего достаточно для решения многих задач
- Однако по эффективности эти оценки не являются наилучшими, их эффективность часто значительно меньше единицы
- ММ часто применяют на практике, т.к. он сводится к простым вычислениям.

Метод максимального правдоподобия (ММП).

В качестве оценки неизвестного параметра θ СВ X выбирается то его значение, при котором полученное значение выборки имеет наибольшую вероятность или плотность вероятности (для дискретных и непрерывных СВ соответственно).

Использование метода сводится к нахождению максимума функции одной или нескольких переменных.

Оценкой наибольшего правдоподобия параметра θ является такое его значение θ^* , при котором функция правдоподобия L имеет максимум.

Алгоритм применения ММП:

- 1) По данной выборке $\{x_n\}$ построить функцию правдоподобия $L(x_1, x_2, \dots, x_n; \theta)$
- 2) Найти натуральный логарифм функции правдоподобия (упростить)
- 3) Найти первую производную $\ln L$
- 4) Решить уравнение правдоподобия $(\ln L)' = 0$
- 5) Найти критические точки θ^* - крит
- 6) Убедиться, что найденная точка θ^* является точкой максимума, воспользовавшись достаточным условием максимума:
 - a. Найти вторую производную $(\ln L)''$
 - b. Если $(\ln L)'' |_{\theta^*} < 0$, то θ^* - точка максимума.

Таким образом найденная точка θ^* будет являться оценкой максимального правдоподобия параметра θ .

Оценки, полученные ММП являются состоятельными и эффективными, но не всегда несмещенными.

Недостатки ММП:

- Решение уравнения правдоподобия бывает трудным
- Нужно знать закон распределения СВ X , что на практике часто бывает невозможным

15. Понятие доверительной вероятности и доверительного интервала. Точность оценки.

Задача интервального оценивания – по данным выборки построить числовой интервал $(\theta_1; \theta_2)$, относительно которого с заранее выбранной вероятностью γ можно сказать, что внутри этого интервала находится оцениваемый параметр.

Интервальной оценкой параметра θ называется числовой интервал $(\theta_1; \theta_2)$, который с заданной вероятностью γ накрывает (содержит) неизвестное значение параметра θ т.е. для которого выполняется условие: $P(\theta_1 < \theta < \theta_2) = \gamma$

Величина γ – есть вероятность события, что параметр θ отличается от своей точечной оценки на величину, не превосходящую некоторого положительного числа δ

$\gamma = P(|\theta - \theta^i| < \delta)$, где δ - точность оценки

Интервал $(\theta_1; \theta_2)$ называется доверительным интервалом или интервальной оценкой, а вероятность γ – доверительной вероятностью, уровнем доверия или надежностью оценки.

Число $\alpha = 1 - \gamma$ – уровень значимости:

- Значение $\alpha(\gamma)$ выбирают заранее
- Выбор зависит от конкретно решаемой задачи
- Обычно выбирают: $\gamma = 0,9; 0,95; 0,99; 0,999$
- Или $\alpha = 0,1; 0,05; 0,01; 0,001$
- Ширина интервала $h = \theta_2 - \theta_1$ существенно зависит от объема выборки n (уменьшается с ростом n) и от значения доверительной вероятности (увеличивается с приближением γ к единице)

16. Основные законы распределения математической статистики.

Распределение Хи-квадрат (Пирсона) с $k=n$ степенями свободы называют распределение суммы квадратов независимых СВ, распределенных по стандартному нормальному закону, т.е.

$Z = Y_1^2 + Y_2^2 + \dots + Y_n^2$, где Y_i – независимы и имеют одно и то же нормальное стандартное распределение, т.е. $Y_i \sim N(0,1)$; $M(Y_i) = 0$, $D(Y_i) = 1$

Число слагаемых называется числом степеней свободы распределения (это параметр формы кривой χ^2 -распределения)

Обозначается $Z \sim \chi^2(n)$

СВ, распределенная по закону χ^2 принимает только неотрицательные значения, т.е. $Z \geq 0$

Кривая χ^2 -распределения одновершинная, асимметричная. С ростом n вершина кривой сдвигается вправо от начала координат и χ^2 - распределение медленно приближается к нормальному.

Числовые характеристики: $M(\chi^2) = n$; $D(\chi^2) = 2n$

Распределение χ^2 используется при интервальном оценивании дисперсии, проверке статистических гипотез согласия, однородности, независимости, при проверке значимости коэффициента корреляции и в других задачах статистического анализа данных.

Распределение Стьюдента или t-распределение с k степенями свободы называется распределение

СВ $Z = \frac{X}{\sqrt{\frac{Y}{k}}}$, где X - СВ, распределенная по стандартному нормальному закону ($X \sim N(0,1)$), а Y

имеет распределение Хи-квадрат с k степенями свободы, т.е. $Y \sim \chi^2(k)$. Обозначение $Z \sim t(k)$. СВ, распределенная по закону Стьюдента принимает любые значения.

Распределение Стьюдента симметрично относительно 0: если $X \sim t(k)$, то и $(-X) \sim t(k)$.

Числовые характеристики: $M(t) = M_0(t) = M_e(t) = 0$; $D(t) = \frac{k}{k-2}$, $k > 2$

Кривая распределения Стьюдента по сравнению с кривой нормального распределения является более пологой. При $k \rightarrow \infty$, t-распределение быстро приближается к нормальному, при $k > 30$ они уже практически неразличимы.

Распределение Стьюдента применяют при интервальном оценивании мат ожидания и других параметров, при проверке гипотез о значениях мат ожиданий, оценке значимости коэф-ов регрессии, гипотез об однородности выборок и тд.

Распределением Фишера с k_1 и k_2 степенями свободы называется распределение СВ

$$Z = \frac{X}{k_1} : \frac{Y}{k_2} = \frac{X * k_2}{Y * k_1}, \text{ где } X \text{ и } Y - \text{ независимые СВ распределенные по закону Пирсона, т.е. } X \sim \chi^2(k_1); Y \sim \chi^2(k_2)$$

Обозначение $Z \sim F(k_1, k_2)$. СВ, распределенная по закону Фишера может принимать только неотрицательные значения, т.е. $Z \geq 0$

Числовые характеристики:

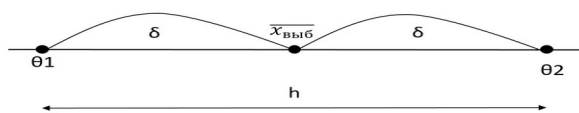
$$M(F) = \frac{k_2}{k_2 - 2}, k_2 > 2$$

$$D(F) = \frac{2 k_2^2 (k_1 + k_2 - 2)}{k_1 (k_2 - 2)^2 (k_2 - 4)}, k_2 > 4$$

С ростом числа степеней свободы распределение Фишера приближается к нормальному. Применяется при проверке гипотез о равенстве дисперсий, о статистической значимости уравнения регрессии.

17. Интервальная оценка генеральной средней (математического ожидания).

В качестве интервальной оценки мат ожидания выступает точечная оценка – выборочная средняя – относительно которой строится симметричный интервал.



Правила построения доверительного интервала зависят от того, известна или неизвестна генеральная дисперсия $\sigma_{ген}^2$

- Доверительный интервал с заданной надежностью γ неизвестного математического ожидания с известной генеральной дисперсией имеет вид:

$$\overline{x}_{выб} - \frac{\sigma_{ген}}{\sqrt{n}} z < a < \overline{x}_{выб} + \frac{\sigma_{ген}}{\sqrt{n}} z$$

Где Z – квантиль стандартного нормального распределения уровня $1 - \alpha/2$

В Excel: =НОРМ.СТ.ОБР((1+ γ)/2)

Число $\delta = \frac{\sigma_{ген}}{\sqrt{n}} z$ – точность оценки мат ожидания. Тогда доверительный интервал для

генеральной средней можно записать в виде: $(\theta_1; \theta_2) = (\overline{x}_{выб} - \delta; \overline{x}_{выб} + \delta)$

Ширина интервала $h = \theta_2 - \theta_1 = 2\delta$

- Доверительный интервал с заданной надежностью γ неизвестного математического ожидания с неизвестной генеральной дисперсией имеет вид:

$\bar{x}_{выб} - \frac{S}{\sqrt{n}} t_{\alpha} < a < \bar{x}_{выб} + \frac{S}{\sqrt{n}} t_{\alpha}$, где t_{α} – квантиль распределения Стьюдента (для двусторонней области) соответствующий $k=n-1$ степеням свободы и уровню значимости α

В Excel: = СТЬЮДЕНТ.ОБР.2Х(α ; $k=n-1$)

Число $\delta = \frac{S}{\sqrt{n}} t_{\alpha}$ – точность оценки мат ожидания.

18. Интервальная оценка генеральной дисперсии и генерального СКО.

Основой интервальной оценки генеральной дисперсии является статистика S^2 (Исправленная выборочная дисперсия или исправленное выборочное СКО – S)

Интервал в отличие от генеральной средней для генеральной дисперсии строится несимметричный.

Доверительный интервал с заданной надежностью γ для генеральной дисперсии имеет вид:

$$\frac{(n-1) * S^2}{\chi_2^2} < \sigma_{ген}^2 < \frac{(n-1) * S^2}{\chi_1^2}$$

Где критические точки определяются по таблице распределения Пирсона:

$$\chi_1^2 = \chi^2 \left(\frac{(1-\gamma)}{2}; n-1 \right)$$

$$\chi_2^2 = \chi^2 \left(\frac{(1+\gamma)}{2}; n-1 \right)$$

Замечание: значения критических точек распределения можно найти с помощью функций Excel:

$$\chi_1^2 = \text{ХИ2.ОБР} \left(\frac{(1-\gamma)}{2}; n-1 \right)$$

$$\chi_2^2 = \text{ХИ2.ОБР} \left(\frac{(1+\gamma)}{2}; n-1 \right)$$

Доверительный интервал для генерального СКО имеет вид:

$$S * \sqrt{\frac{n-1}{\chi_2^2}} < \sigma_{ген} < S * \sqrt{\frac{n-1}{\chi_1^2}}$$

19. Интервальная оценка генеральной доли.

Доверительный интервал для генеральной доли с заданной надежностью γ строится симметрично относительно выборочной доли и имеет вид: $\omega - \delta < p < \omega + \delta$

Где $\omega = \frac{m}{n}$ – выборочная доля

$$\delta = z * \sqrt{\frac{\omega(\omega-1)}{n}}$$

точность оценки генеральной доли

Z - квантиль стандартного нормального распределения уровня $1 - \frac{\alpha}{2}$

$$Z = \text{НОРМ.СТ.ОБР} \left(\frac{(1+\gamma)}{2} \right)$$

20. Понятие статистической гипотезы. Параметрические и непараметрические гипотезы.

Статистическая гипотеза – это некоторое высказывание относительно генеральной совокупности, проверяемое по выборочным данным.

Статистическая гипотеза – это любое предположение о виде неизвестного распределения или о параметрах известного закона распределения.

Статистические гипотезы:

- Параметрические
 - Простые
 - Сложные
- Непараметрические

Параметрические гипотезы – это утверждения о значении параметров распределения известного вида:

- Простые – утверждения о том, что значение неизвестного параметра генеральной совокупности равно одному заданному числу
- Сложные – гипотезы, не являющиеся простыми

Непараметрические гипотезы – это утверждения о виде неизвестного распределения

21. Нулевая и альтернативная гипотеза. Типы альтернативных гипотез.

Выдвинутая гипотеза называется нулевой (основной). H_0 – нулевая гипотеза. По отношению к основной гипотезе можно сформулировать альтернативную (конкурирующую) противоречащую ей, т.е. гипотезу противоположную H_0 . H_1 – альтернативная гипотеза.

H_0 и H_1 – это два предположения, из которых в результате статистической проверки должно быть выбрано только одно.

Проверить статистическую гипотезу значит установить, согласуются ли данные, полученные из выборки с этой гипотезой.

Статистическими методами гипотезу можно только опровергнуть или не опровергнуть, но не доказать!!!

Примеры параметрических гипотез:

- H_0 : мат ожидание СВ X равно 2 (нулевая, основная)
- H_1 : $a \neq 2$ – двусторонняя сложная гипотеза
- H_1 : $a > 2$ – правосторонняя сложная гипотеза
- H_1 : $a < 2$ – левосторонняя сложная гипотеза
- H_1 : $a = 1,9$ – простая левосторонняя гипотеза
- H_1 : $a = 2,1$ – простая правосторонняя гипотеза

22. Задача проверки статистических гипотез. Понятие статистического критерия.

Сопоставление высказанной гипотезы о генеральной совокупности с имеющимися выборочными данными, сопровождаемое количественной оценкой степени достоверности получаемого вывода и осуществляемое с помощью того или иного статистического критерия называется проверкой статистических гипотез.

Статистический критерий K (или стат тест) – это правило (формула), по которому определяется мера расхождения результатов выборочного наблюдения с высказанной гипотезой H_0 , т.е. согласно которому принимается решение, принимать или отклонять нулевую гипотезу.

Основу критерия представляет специально составленная выборочная характеристика (СВ или статистика) $\theta_n = \theta(X_1, X_2, \dots, X_n)$ - точечное или приближенное распределение которой неизвестно.

23. Наблюдаемое и критическое значения статистического критерия.

Значение критерия, рассчитываемое по специальным правилам (по формулам) на основании выборочных данных называют наблюдаемым значением критерия ($K_{\text{набл}}$)

Множество значений критерия K :

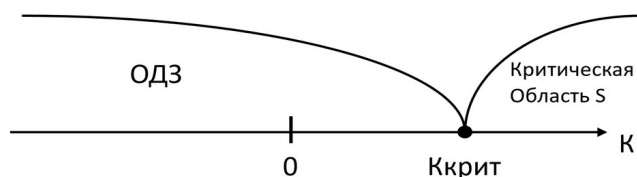
- Область допустимых значений (область принятия нулевой гипотезы) \bar{S} – совокупность всех значений критерия K , при которых нулевая гипотеза H_0 не отклоняется
- Критическая область S – совокупность всех значений критерия K , при которых нулевая гипотеза отклоняется в пользу конкурирующей H_1 .

Значения критерия, разделяющие области S и \bar{S} , определяемые на заданном уровне значимости α по таблицам распределения (инструментальными средствами) СВ K , выбранной в качестве критерия, называют критическими точками ($K_{\text{крит}}$)

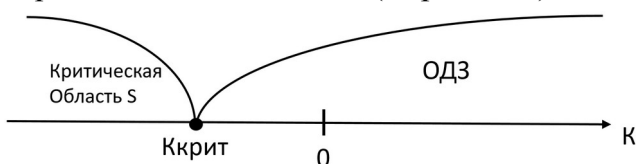
24. Критическая область и ее типы. Область принятия гипотезы.

Вид критической области зависит от того, какая гипотеза выдвинута в качестве альтернативной:

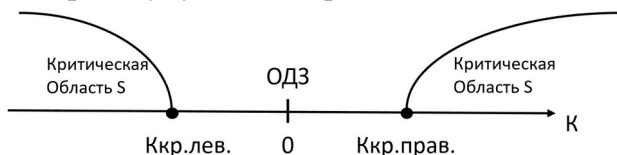
- Если конкурирующая гипотеза правосторонняя ($H_1: a > 2$), то и критическая область S правосторонняя. В этом случае будет одна критическая точка, которая принимает положительные значения ($K_{\text{крит.прав.}} > 0$)



- Если конкурирующая гипотеза левосторонняя ($H_1: a < 2$), то и критическая область S левосторонняя. В этом случае будет одна критическая точка, принимающая отрицательные значения. ($K_{\text{кр.лев.}} < 0$)



- Если конкурирующая область двусторонняя ($H_1: a \neq 2$), то и критическая область S будет двусторонняя. В этом случае будет 2 критических точки ($K_{\text{кр.лев.}} < 0$ и $K_{\text{кр.прав.}} > 0$), которые будут симметричны относительно нуля.



25. Основной принцип проверки статистических гипотез.

Основной принцип проверки статистических гипотез:

- Если наблюдаемое значение критерия ($K_{\text{набл}}$) принадлежит критической области S , то нулевая гипотеза отклоняется в пользу конкурирующей гипотезы H_1
- Если наблюдаемое значение критерия ($K_{\text{набл}}$) принадлежит ОДЗ, то нулевую гипотезу H_0 нельзя отклонять.

26. Общая схема проверки статистических гипотез.

- 1) Располагая выборкой x_1, x_2, \dots, x_n сформулировать основную (H_0) и альтернативную (H_1) гипотезы
- 2) Выбрать уровень значимости α для проведения проверки
- 3) По виду H_0 выбрать статистический критерий для ее проверки
- 4) На основании данных выборки по специальному алгоритму (формуле) найти наблюдаемое значение критерия $K_{набл}$
- 5) По таблицам (формулам Excel) распределения СВ K , выбранной в качестве статистического критерия, найти критическое значение (критическую точку или точки)
- 6) Исходя из типа конкурирующей гипотезы H_1 определить тип критической области S
- 7) Посмотреть, какой области (S или \bar{S}) принадлежит наблюдаемое значение критерия
- 8) Сделать вывод о принятии или отклонении нулевой гипотезы.

Замечание: даже в том случае, если H_0 нельзя отклонить, это еще не значит, что данное предположение о генеральной совокупности является единственным верным, просто ему не противоречат имеющиеся выборочные данные. Таким свойством могут обладать и другие гипотезы.

27. Ошибки первого и второго рода. Мощность критерия.

При проверке гипотезы могут быть приняты неправильные решения, т.е. могут быть допущены ошибки первого и второго рода.

Ошибкой первого рода называется ошибка, возникающая, когда нулевая гипотеза отклоняется в то время как в действительности в генеральной совокупности она является справедливой. Вероятность совершить ошибку первого рода называется уровнем значимости (или размером критерия) и обозначается α .

Ошибкой второго рода называется ошибка, возникающая, когда H_0 принимается, в то время как на самом деле в генеральной совокупности она является ошибочной, а справедлива альтернативная H_1 .

Вероятность допустить ошибку второго рода называется β

Вероятность $1 - \beta$, т.е. вероятность не допустить ошибку второго рода называют мощностью критерия (или функцией мощности)

Выбор статистического критерия и критической области осуществляют таким образом, чтобы мощность критерия была максимальной.

Критерий называют наиболее мощным, если из всех критериев с заданным уровнем значимости α , он имеет наибольшую мощность (наименьшее значение β)

Риски при проверке гипотез

Нулевая гипотеза H_0	Результаты решения относительно нулевой гипотезы	
	Отклонена	Принята
Верна	Ошибка первого рода $P(H_1/H_0)=\alpha$ (уровень значимости)	Правильное решение $P(H_0/H_0)=1-\alpha$ (надежность)
Неверна	Правильное решение $P(H_1/H_1)=1-\beta$ (мощность критерия)	Ошибка второго рода $P(H_0/H_1)=\beta$

28. Наблюдаемый уровень значимости (p-значение/p-value)

Наблюдаемым уровнем значимости (или р-значением или р-value) гипотезы H_0 , проверяемой по выборке x_1, x_2, \dots, x_n с помощью критерия $K(x_1, x_2, \dots, x_n)$ называется наименьшее значение α , при котором нулевая гипотеза отклоняется

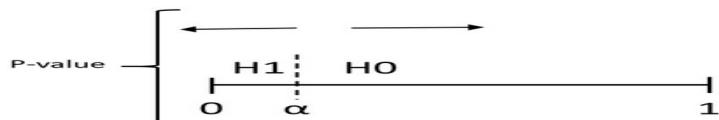
$$p(x_1, x_2, \dots, x_n) = \min \{ \alpha : K(x_1, x_2, \dots, x_n) \in S \}$$

Ещё одна интерпретация р-значения – это вероятность, с которой (при условии истинности H_0) могла бы реализоваться полученная выборка или любая другая выборка с еще менее вероятным наблюдаемым значением статистики K .

СВ $p(x_1, x_2, \dots, x_n)$ имеет равномерное распределение.

- Если $p(x_1, x_2, \dots, x_n) < \alpha$, то есть основания отклонять H_0
- Если $p(x_1, x_2, \dots, x_n) \geq \alpha$, то нет оснований отклонять H_0

Поэтому на практике, чем меньше р-значение, тем меньше вероятность ошибиться, отклонив нулевую гипотезу и тем выше уверенность в том, что необходимо отвергнуть H_0 .



29. Гипотезы о равенстве средних 2-х нормально распределенных генеральных совокупностей.

Случай 1: генеральные дисперсии известны (z-тест)

- 1) Найдем наблюдаемое значение критерия. Найдем точечные оценки мат ожидания (генеральные средние). При достаточно больших объемах выборок (≥ 30) выборочные средние имеют приближенно нормальный закон распределения, поэтому при выполнении H_0 статистика Z имеет стандартный нормальный закон распределения, т.е. формула (1):

$$Z = \frac{\bar{x}_{выб} - \bar{y}_{выб}}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}} \sim N(0; 1)$$

По формуле (1) находим наблюдаемое значение критерия.

- 2) Найдем критические точки и критическую область S . Критическая область зависит от типа альтернативной гипотезы H_1 .
 - a. H_1 – двусторон. S -двусторон. Найдем 2 симметричные критические точки:
 $Z_{кр.пр.}$ – квантиль станд норм расп. Уровня $1-\alpha/2$ (=НОРМ.СТ.ОБР($1-\alpha/2$))
 $Z_{кр.лев.}$ – квантиль станд норм расп уровня $\alpha/2$ (= $-Z_{кр.пр.}$)
 $S = (-\infty; Z_{кр.лев.}) \cup (Z_{кр.пр.}; +\infty)$
 - b. H_1 - правосторон S -правосторон. Найдем одну критическую точку (правую):
 $Z_{кр.пр.}$ – квантиль станд норм расп уровня $1-\alpha$ (=НОРМ.СТ.ОБР($1-\alpha$))
 $S = (Z_{кр.пр.}; +\infty)$

- 3) Принятие решения относительно нулевой гипотезы: Если Z принадлежит критической области, то H_0 отклоняется в пользу H_1 . Иначе нет оснований отклонять H_0 .

Замечание 1: Если в условии даны несгруппированные ряды, то Z-тест в Excel можно реализовать с помощью инструмента в пакете анализа «Двухвыборочный Z-тест для средних».

Замечание 2: Этот инструмент дает значение р-value только для односторонней H_1 , поэтому если H_1 -двусторон, то полученное р-value надо умножить на 2.

Случай 2: генеральные дисперсии неизвестны (t-тест – критерий Стьюдента)

1) Найдем несмещенные точечные оценки генеральных средних и генеральных СКО.

$$a_x^{\dot{}} = \overline{x_{выб}}$$

$$a_y^{\dot{}} = \overline{y_{выб}}$$

$$\sigma_x^{\dot{}} = S_x$$

$$\sigma_y^{\dot{}} = S_y$$

При условии выполнения H_0 , статистика t имеет распределение Стьюдента с $k=n_x+n_y-2$ степенями свободы, т.е. формула (2):

$$t = \frac{\overline{x_{выб}} - \overline{y_{выб}}}{S * \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} \quad t(k)$$

$$S = \sqrt{\frac{(n_x - 1) * S_x^2 + (n_y - 1) * S_y^2}{n_x + n_y - 2}} \text{ – объединенная оценка СКО}$$

По формуле (2) находим t набл

2) Найдем критические точки и критическую область:

a. H_1 – двусторон. S- двусторон.

Найдем 2 симметричные критические точки.

$t_{кр.пр.}$ - квантиль распределения Стьюдента уровня $1-\alpha/2$ (=СТЮДЕНТ.ОБР(1- $\alpha/2$;k)
=СТЮДЕНТ.ОБР.2X(α ;k))

$t_{кр.лев.}$ - квантиль распределения Стьюдента уровня $\alpha/2 = -t_{кр.пр.}$

$$S = (-\infty; t_{кр.лев.}) \cup (t_{кр.пр.}; +\infty)$$

b. H_1 -правосторон S-правосторон

Найдем одну критическую точку (правую)

$t_{кр.пр.}$ - квантиль распределения Стьюдента уровня $1-\alpha$ (=СТЮДЕНТ.ОБР(1- α ;k)
=СТЮДЕНТ.ОБР.ПХ(α ;k))

$$S = (t_{кр.пр.}; +\infty)$$

3) Принятие решения относительно нулевой гипотезы. Если наблюдаемое значение t принадлежит критической области, то H_0 отклоняется в пользу H_1 на уровне значимости альфа, иначе нет оснований отклонять H_0 .

Замечание 1: Описанный способ проверки с помощью объединенного двухвыборочного t-теста можно проводить только когда генеральные дисперсии совокупностей X и Y можно считать равными, поэтому перед проверкой H_0 нужно сначала проверить гипотезу о равенстве генеральных дисперсий с помощью критерия Фишера (F-теста).

Замечание 2: Если в условии даны несгруппированные выборки, то проверку можно осуществить с помощью инструментов Excel:

- Если Генеральные дисперсии можно считать равными, то используем «двухвыборочный t-тест с одинаковыми дисперсиями
- Если генеральные дисперсии нельзя считать равными, то используем «двухвыборочный t-тест с различными дисперсиями» (!!!р-значение в этом случае считается неверно)

Замечание 3: Если не проводить проверку гипотезы о равенстве генеральных дисперсий (неизвестно равны они или нет), то t-тест можно реализовать с помощью функций Excel.

P-value:

- =СТЮДЕНТ.ТЕСТ(МассивX ; Массив Y; Хвосты (1, если H1 – односторонняя, 2, если двусторонняя); Тип (3))
- =ТТЕСТ(аргументы те же)

Если p-значение меньше альфа, то H0 отклоняется в пользу H1, иначе нет оснований отклонять H0.

30. Гипотеза о равенстве дисперсий 2-х нормально распределенных генеральных совокупностей.

1) Найдем наблюдаемое значение критерия:

Найдем точечные несмещенные оценки генеральных дисперсий (S^2)

Если H0 выполняется, то F-статистика имеет распределение Фишера с k_1, k_2 степенями свободы, т.е. по формуле (3):

$$F = \frac{S_{\text{большая}}^2}{S_{\text{меньшая}}^2} \sim F(k_1, k_2)$$

Где $k_1 = n_1 - 1$ (n_1 - объем выборки с большей исправленной выб дисперсией), $k_2 = n_2 - 1$ (n_2 – объем выборки с меньшей исп выб дисп)

По формуле (3) найдем наблюдаемое значение критерия.

2) Найдем критические точки и критическую область

- Если H1 двусторон, то S- двусторон.

$F \geq 0$, найдем 2 критические точки.

$F_{\text{кр.лев.}}$ – квантиль распределения Фишера уровня $\alpha/2$ ($=F.OБР(\alpha/2; k_1; k_2) < 1$)

$F_{\text{кр.пр}}$ – квантиль распределения Фишера уровня $1 - \alpha/2$ ($=F.OБР(1 - \alpha/2; k_1; k_2) > 1$)

$S = (0; F_{\text{кр.лев}}) \cup (F_{\text{кр.пр.}}; +\infty)$

- Если H1 правосторон, то S-правосторон

$F \geq 0$, найдем $F_{\text{кр пр}}$

$F_{\text{кр.пр}}$ – квантиль распределения Фишера уровня $1 - \alpha$ ($=F.OБР(1 - \alpha; k_1; k_2)$;

$=F.OБР.ПХ(\alpha; k_1; k_2)$)

$S = (F_{\text{кр.пр}}; +\infty)$

3) Принятие решения относительно H0: Если наблюдаемое значение принадлежит критической области, то H0 отклоняется в пользу H1, иначе нет оснований отклонять H0.

Замечание: По несгруппированным выборкам F- тест можно провести следующими способами:

- С помощью инструмента надстройки анализа данных «двухвыборочный F-тест для дисперсий»:
 - Первый массив нужно ввести массив с наибольшей исп выб дисп
 - Если H1 односторонняя, то в строку «Альфа» вводим значение α , если двусторонняя, то $\alpha/2$
 - Этот инструмент считает p-значение для односторонней H1, если H1- двусторонняя, то p-значение надо умножить на 2
- С помощью функции =F.ТЕСТ(Массив1;Массив2) – считает p-значение только для двусторонней H1, если H1 односторонняя, то значение надо разделить на 2.

31. Гипотеза о равенстве генеральной долей признака в двух совокупностях.

Сравнение вероятностей успеха в двух сериях испытания Бернулли (z-тест)

Имеется 2 выборки, распределенные по биномиальному закону с параметрами $(n_i; p_i)$, где p_i – вероятность успеха в одном испытании Бернулли.

Проверка:

1) Найдем наблюдаемое значение критерия.

По выборкам $A:n_1;k_1$, где k_1 – число элементов с признаком и $B:n_2;k_2$. n_1 и $n_2 \geq 30$.

Найдем несмещенные точечные оценки генеральных долей, т.е. выборочные доли:

$$p_1^{\hat{}} = \omega_1 = \frac{k_1}{n_1}$$

$$p_2^{\hat{}} = \omega_2 = \frac{k_2}{n_2}$$

При достаточно больших объемах выборок выборочные доли имеют приближенно нормальный закон распределения. При условии, что H_0 выполняется Z-статистика имеет стандартное нормальное распределение, т.е. по формуле (4):

$$Z = \omega_1 - \omega_2 - \frac{1}{2} \sqrt{\omega_1 \omega_2}$$

$$\omega = \frac{k_1 + k_2}{n_1 + n_2} \text{ – общая доля из двух выборок}$$

По формуле 4 находим Z набл.

2) Найдем критические точки и критическую область

- Если H_1 – двусторон, S – двусторон,
Найдем 2 симметричные критические точки
Zкр.пр. – квантиль станд норм расп уровня $1-\alpha/2$ (=НОРМ.СТ.ОБР($1-\alpha/2$))
Zкр.лев. – квантиль станд норм расп уровня $\alpha/2$ (-Zкр.пр.)
 $S = (-\infty; Z_{кр.лев}) \cup (Z_{кр.пр.}; +\infty)$
- Если H_1 – правосторон, S – правосторон
Найдем Zкр пр
Zкр. Пр. – квантиль станд норм рас уровня $1-\alpha$ (=НОРМ.СТ.ОБР($1-\alpha$))
 $S = (Z_{кр.пр.}; +\infty)$

3) Принятие решения

Замечание: Z-тест можно использовать только при больших объемах выборки, т.е. когда выборки по объему репрезентативны (представительны):

- $n_1 * \omega_1 \geq 5$
- $n_1 * (1 - \omega_1) \geq 5$
- $n_2 * \omega_2 \geq 5$
- $n_2 * (1 - \omega_2) \geq 5$

Должны выполняться все 4 условия

32. Гипотезы о числовых значениях параметров распределения.

33. Суть и схема применения критерия согласия Пирсона.