

Доверительные интервалы для доли и дисперсии

План

- Доверительный интервал для доли
- Доверительный интервал для дисперсии
- Статистический бутстреппинг

Доверительный интервал для доли

- Описание проблемы
- ДИ
- Алгоритм
- Пример

Оценка доли признака

Задача состоит в построении доверительной оценки для генеральной доли, если известно значение выборочной доли.

Пример. Среди 500 резюме кандидатов на работу няни оказалось 60 принадлежащих мужчинам. Если считать, что выборка репрезентативна, то требуется построить 90%-ый доверительный интервал для фактической доли мужчин, устраивающихся на работу нянями.

Оценка доли признака

	Генеральная совокупность	Выборочная совокупность
Общее число объектов	N	n
Частота	pN	m
Доля признака	p Параметр	$\hat{p} = \frac{m}{n}$ Оценка

Формальное описание проблемы

Цель. Оценить долю признака в генеральной совокупности.

Что мы имеем. Имеем случайную выборку объема n из генеральной совокупности. По выборке вычислена доля признака. Выполнены условия $np \geq 5$ и $n(1 - p) \geq 5$.

Требуется. Построить доверительный интервал для доли:

$$p - E < \hat{p} < p + E$$

Доверительный интервал для доли

Доля значений признака в генеральной совокупности с надежностью $1 - \alpha/2$ находится в доверительном интервале:

$$\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}} < p < \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$$

Обозначение: $\hat{q} = 1 - \hat{p}$

Последовательность действий

- Шаг 1.** По выборке вычислить долю признака.
- Шаг 2.** По таблице нормального распределения найти z-значение для доверительной вероятности $1 - \alpha$.
- Шаг 3.** Вычислить точность интервальной оценки по формуле:

$$E = z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$$

Последовательность действий

Шаг 4. Подставить полученные значения в формулу для доверительного интервала:

$$p - E < p < p + E$$

Шаг 5. Написать ответ.

Пример. Выборы мэра

В ходе проведенного опроса 829 жителей города выяснилось, что 417 опрошенных (51,5%) предполагают поддержать на предстоящих выборах кандидатуру действующего мера.

Местная многотиражка поспешила заявить, что более половины жителей города поддерживают перевыборы действующего мера на следующий срок.

Построить доверительный интервал для доли генеральной совокупности и проверить утверждение корреспондента.

Решение

Шаг 1. По условию, доля признака в выборке составила:

$$\hat{p} = 0,515$$

$$\hat{q} = 0,485$$

$$n = 829$$

Шаг 2. Для доверительной вероятности $1 - \alpha = 0,95$ по таблице нормального закона находим z-значение:

$$z_{\alpha/2} = 1,96$$

Решение

Шаг 3. Вычисляем точность интервальной оценки:

$$E = 1,96 \cdot \sqrt{\frac{0,515 \cdot 0,485}{829}} = 0,034$$

Шаг 4. Подставляем полученные значения в формулу для доверительного интервала:

$$0,515 - 0,034 < p < 0,515 + 0,034$$

Шаг 5. Ответ:

$$48,1\% < p < 54,9\%$$

Пример. Мужчины-няни

Среди 500 резюме кандидатов на работу няни оказалось 60 принадлежащих мужчинам.

Найти 90%-ый доверительный интервал для фактической доли мужчин, устраивающихся работать нянями.

Решение

Шаг 1. По условию, доля признака в выборке составила:

$$\hat{p} = 0,12$$

$$\hat{q} = 0,88$$

$$n = 500$$

Шаг 2. Для доверительной вероятности $1 - \alpha = 0,90$ по таблице нормального закона находим z-значение:

$$z_{\alpha/2} = 1,65$$

Решение

Шаг 3. Вычисляем точность интервальной оценки:

$$E = 1,65 \cdot \sqrt{\frac{0,12 \cdot 0,88}{500}} = 0,024$$

Шаг 4. Подставляем полученные значения в формулу для доверительного интервала:

$$0,12 - 0,024 < p < 0,12 + 0,024$$

Шаг 5. Ответ:

$$9,6\% < p < 14,4\%$$

Объем выборки для оценки доли

Минимальный объем выборки, требуемый для интервального оценивания генеральной доли, находится по формуле:

$$n = \hat{p}\hat{q} \cdot \left(\frac{z_{\alpha/2}}{E} \right)^2$$

При необходимости следует округлить n , чтобы получить целое число.

Важное замечание

Если оценка для доли неизвестна, минимальный объем находят по формуле:

$$n = 0,25 \cdot \left(\frac{z_{\alpha/2}}{E} \right)^2$$

Пример. У кого есть дома компьютер?

Исследователь хочет с 95%-ой вероятностью оценить количество людей, у которых дома имеется персональный компьютер.

По данным предыдущего исследования у 40% опрошенных есть дома компьютер.

Исследователь не хочет ошибиться больше, чем на 2% по сравнению с генеральной долей.

Найти минимальный размер выборки.

Решение

Поскольку $1 - \alpha = 0,95$, то z-значение равно 1,96. $E = 0,02$.

Подставляем в формулу и вычисляем:

$$\begin{aligned} n &= \hat{p}\hat{q} \cdot \left(\frac{z_{\alpha/2}}{E} \right)^2 = \\ &= 0,40 \cdot 0,60 \cdot \left(\frac{1,96}{0,02} \right)^2 = 2304,96 \end{aligned}$$

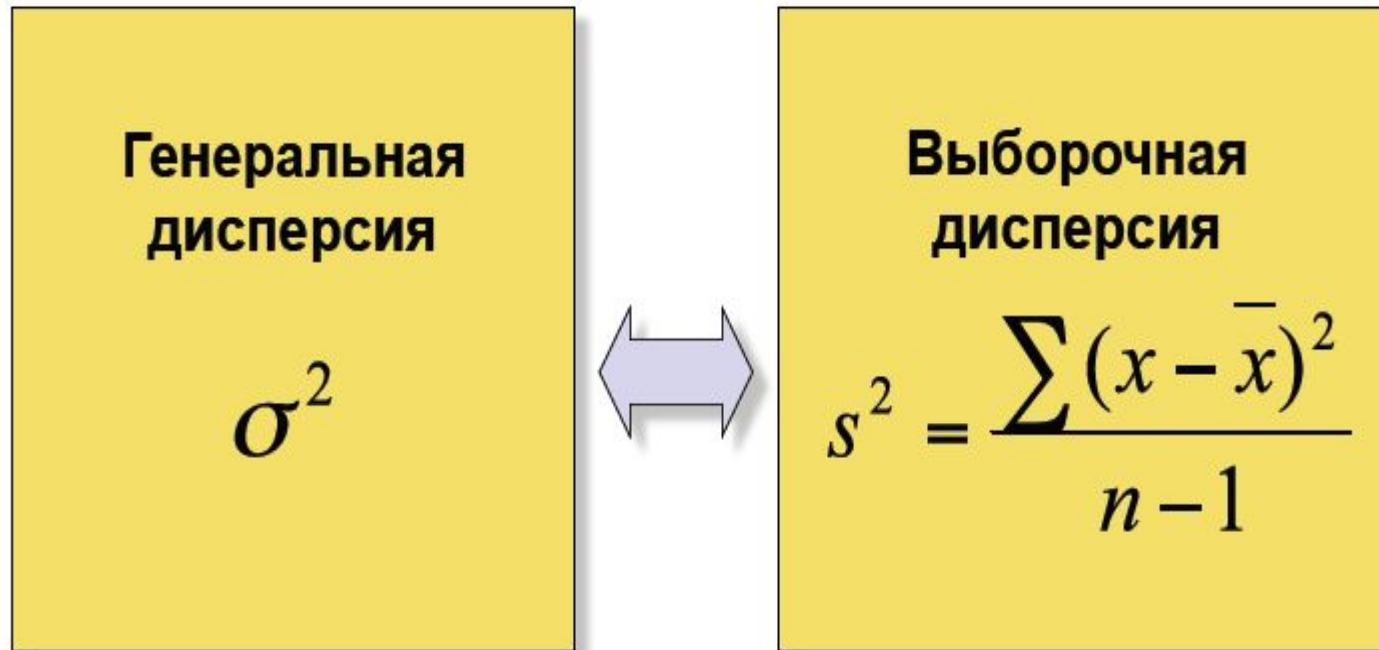
Ответ. Нужно опросить 2305 людей.

Доверительный интервал для дисперсии

- Описание проблемы
- Доверительный интервал
- Алгоритм
- Пример

Оценка для генеральной дисперсии

Задача состоит в построении интервальной оценки генеральной дисперсии на основе выборочной дисперсии.



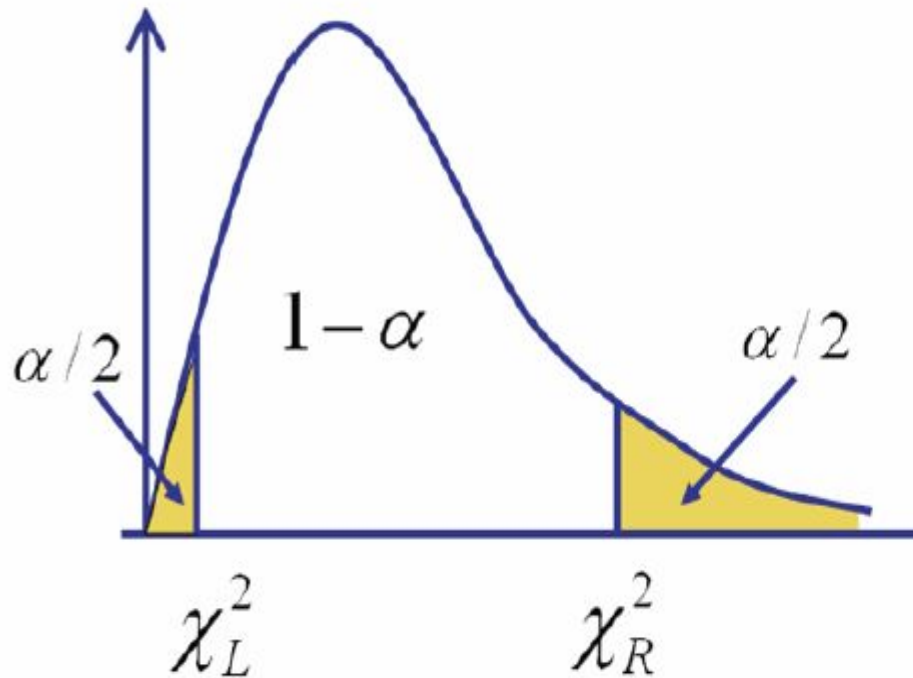
Доверительный интервал для дисперсии

Доверительный интервал для дисперсии находится по формуле:

$$\frac{(n-1)s^2}{\chi_R^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_L^2}$$

Значения хи-квадрат

Значения χ^2_L и χ^2_R находятся по таблицам хи-квадрат распределения, исходя из следующих условий:



$$P(\chi^2 > \chi_L) = 1 - \frac{\alpha}{2}$$

$$P(\chi^2 > \chi_R) = \frac{\alpha}{2}$$

Оценка стандартного отклонения

Доверительный интервал для стандартного отклонения находится по следующей формуле:

$$\frac{\sqrt{n-1} \cdot s}{\chi_R} < \sigma < \frac{\sqrt{n-1} \cdot s}{\chi_L}$$

Последовательность действий

Шаг 1. По выборке вычислить дисперсию.

Шаг 2. По таблице найти два хи-квадрат значения χ^2_L и χ^2_R для доверительной вероятности $1 - \alpha$ и числа степеней свободы $df = n - 1$.

Шаг 3. Подставить полученные значения в формулу:

$$\frac{(n-1)s^2}{\chi_R^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_L^2}$$

Шаг 4. Написать ответ.

Пример. Оценка для дисперсии

Из нормально распределенной генеральной совокупности сделана выборка из 10 элементов. Выборочная дисперсия оказалась равна 28,2.

Требуется оценить дисперсию генеральной совокупности (построить доверительный интервал).

Доверительную вероятность выберем на уровне 90%.

Последовательность действий

Шаг 1. По выборке объема 10 вычислена дисперсия 28,2.

Шаг 2. Число степеней свободы $df = n - 1 = 9$. Поскольку доверительная вероятность равна 90%, по таблице хи-квадрат распределения находим значения $\chi^2_L = 3,325$ и $\chi^2_R = 16,919$.

Шаг 3. Подставим полученные значения в формулу:

$$\frac{(10 - 1) \cdot 28,2}{16,919} < \sigma^2 < \frac{(10 - 1) \cdot 28,2}{3,325}$$

Шаг 4. Ответ:

$$15,0 < \sigma^2 < 76,3$$

Оценка для стандартного отклонения

Если дисперсия находится в доверительном интервале:

$$15,0 < \sigma^2 < 76,3$$

То стандартное отклонение можно оценить следующим образом:

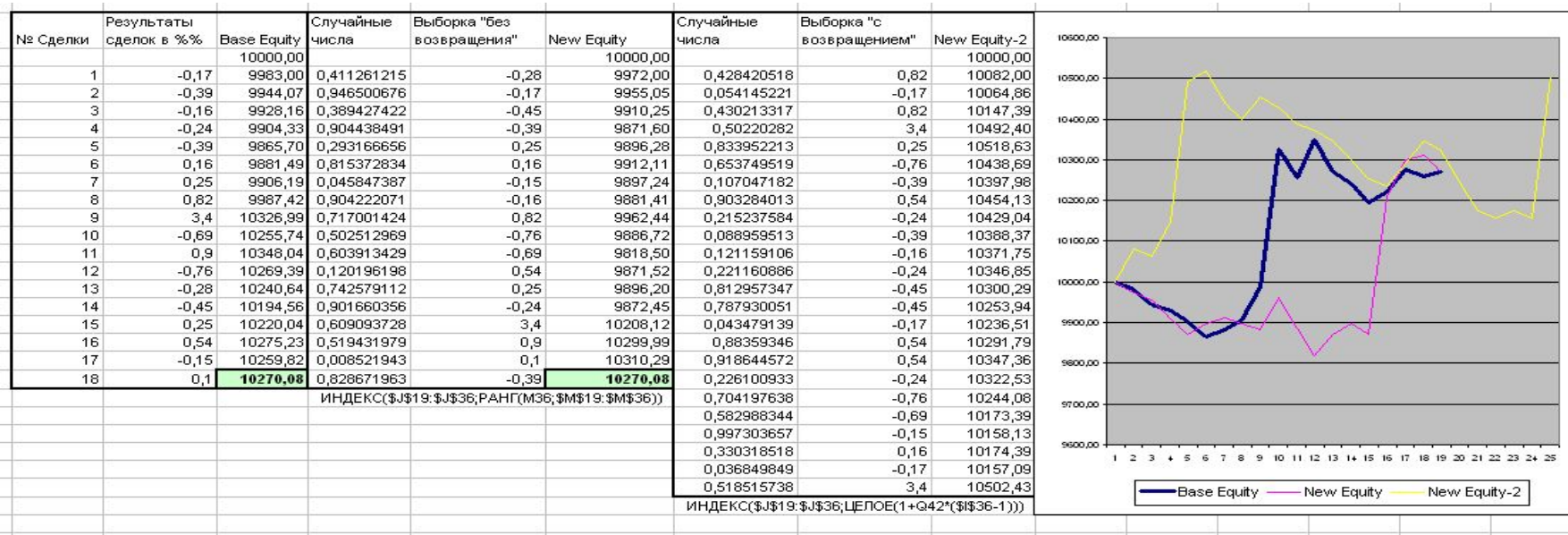
$$\sqrt{15,0} < \sigma < \sqrt{76,3}$$

$$3,87 < \sigma < 8,73$$

Статистический бутстреппинг

- компьютерный метод определения статистик вероятностных распределений. Основан на многократном генерировании выборок методом Монте-Карло на базе данных обучающей выборки

Позволяет просто и быстро оценивать самые разные статистики (доверительные интервалы, дисперсию, корреляцию и так



Методы размножения выборок (бутстреп-методы)

- предложен в 1977 г. Б.Эфроном из Станфордского университета (США)

- "bootstrap" - кожаные петельки на задниках ботинок



- "lift himself by his bootstraps" - "вытащить себя из болота за ушки на задниках ботинок" , "выбиться в люди благодаря собственным усилиям"

Принцип статистического бутстрэппинга

- имитировать многократное получение выборки из генеральной совокупности, используя данные из имеющейся у нас выборки.

Предположим, что мы исследуем высоту людей во всем мире. Мы не можем измерить всех людей, а вместо этого выбираем лишь малую часть. Пусть в нашей выборке N людей. Мы можем посчитать среднее значение. Но для того, чтобы рассуждать о доверительном интервале роста населения, нам нужно некоторое представление о вариабельности среднего.

Используя наши исходные данные о росте N различных людей, составляем новую выборку, также размера N . Это новая выборка взята из исходной случайным образом так, что мы каждый раз случайным образом выбираем из N имеющихся значений). У такой выборки будет другое среднее.

Сделав такую выборку много раз (возможно, 1000 или 10000 раз), каждый раз вычисляя среднее, мы получаем гистограмму распределения, которая может ответить на вопросы о доверительном интервале.