



Прикладная статистика и анализ данных

Съезд VII

Дисперсионный анализ I





Типы задач Д.А.

1. Независимые выборки

Две группы пациентов. Одним дают одно лекарство, другим — другое. Верно ли, что первое лекарство эффективнее?

2. Связные выборки

Пациент проходит испытание, принимает средство, затем снова проходит испытание. Отличается ли эффект?

- ▶ Методы для задач 2 типа можно использовать для задач 1 типа. При этом теряется важная информация.
- ▶ Методы для задач 1 типа *нельзя* использовать для задач 2 типа.



Независимые выборки

Человек	Препарат	Изменение температуры
Петя	Апотивадом	-0.9
Вася	Апотивадом	-0.6
Катя	Апотивадом	-1.0
Миша	Апотивадом	-0.3
Ира	Волымикер	-2.6
Света	Волымикер	-1.9
Коля	Волымикер	-0.7

Значимо ли отличается эффект от приема препаратов?



Связные выборки

Каждый человек применяет один и тот же препарат.

Человек	Температура до	Температура после
Петя	38.2	37.6
Вася	37.6	38.0
Катя	38.5	37.1
Миша	38.0	36.9
Ира	37.9	37.1
Света	39.4	37.3

Есть ли эффект от приема препарата?



Другие вопросы на практике

1. Значимо ли отличаются решения в топ-10 в соревновании на Kaggle?
2. На какой дизайн кнопки клиент кликнет с большей вероятностью?
3. Увеличился ли средний чек корзины покупателей после внедрения нового блока рекомендаций?
4. В чем причина оттока клиентов к конкурентам?
5. Отличаются ли гены по степени экспрессии?
6. многие другие...



Бернуллиевские выборки



Независимые выборки

$X_1, \dots, X_n \sim \text{Bern}(p)$ и $Y_1, \dots, Y_m \sim \text{Bern}(q)$.

$H_0: p = q$ vs. $H_1: p \{<, \neq, >\} q$

$\hat{p}_1 = \bar{X} \stackrel{d}{\approx} \mathcal{N}\left(p_1, \frac{p_1(1-p_1)}{n}\right)$ и $\hat{p}_2 = \bar{Y} \stackrel{d}{\approx} \mathcal{N}\left(p_2, \frac{p_2(1-p_2)}{m}\right)$ — ОМП

При справедливости H_0 : $W(X, Y) = \frac{\hat{p}_1 - \hat{p}_2}{\hat{\sigma}} \stackrel{d}{\approx} \mathcal{N}(0, 1)$,

где $\hat{\sigma} = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n} + \frac{\hat{p}_2(1-\hat{p}_2)}{m}}$.

Для $H_1: p > q$ критерий Вальда $S = \{W(x, y) > z_{1-\alpha}\}$.

Дов. интервал для $p_1 - p_2$ равен $C = (\hat{p}_1 - \hat{p}_2 - z_{1-\alpha}\hat{\sigma}, 1)$.

H_0 отвергается $\iff 0 \notin C$



Пример (влияние нового препарата на выздоровление)

Испытуемые делятся случайно на две группы:

1. *Исследуемая группа* — принимает новый препарат;

$X = (X_1, \dots, X_n) \sim \text{Bern}(p_1)$ — результаты лечения.

2. *Контрольная группа* — принимает плацебо;

$Y = (Y_1, \dots, Y_m) \sim \text{Bern}(p_2)$ — результаты лечения.

$H_0: p_1 = p_2$ — отсутствие эффекта

$H_1: p_1 > p_2$ — эффект присутствует



Пример (влияние нового препарата на выздоровление)

1. 1 группа: $n = 30$ человек, 27 выздоровело $\implies \hat{p}_1 = 0.9$

2 группа: $m = 30$ человек, 21 выздоровело $\implies \hat{p}_2 = 0.7$

$W(x, y) \approx 2$, $pvalue = 0.0228$, дов. интервал $(0.036, 1)$

2. 1 группа: $n = 30$ человек, 27 выздоровело $\implies \hat{p}_1 = 0.9$

2 группа: $m = 30$ человек, 15 выздоровело $\implies \hat{p}_2 = 0.5$

$W(x, y) \approx 3.76$, $pvalue = 0.00008$, дов. интервал $(0.225, 1)$

3. 1 группа: $n = 30$ человек, 27 выздоровело $\implies \hat{p}_1 = 0.9$

2 группа: $m = 10$ человек, 7 выздоровело $\implies \hat{p}_2 = 0.7$

$W(x, y) \approx 1.54$, $pvalue = 0.0618$, дов. интервал $(-0.017, 1)$

Связные выборки

$$X_1, \dots, X_n \sim \text{Bern}(p)$$

$$Y_1, \dots, Y_n \sim \text{Bern}(q).$$

$$H_0: p = q \text{ vs. } H_1: p \{<, \neq, >\} q$$

	$Y_i = 1$	$Y_i = 0$
$X_i = 1$	e	f
$X_i = 0$	g	h

Статистика **критерия**:

$$Z(X, Y) = \frac{\hat{p} - \hat{q}}{\sqrt{\frac{f+g}{n^2} + \frac{(f-g)^2}{n^3}}} = \frac{f - g}{\sqrt{f + g + \frac{(f-g)^2}{n}}} \xrightarrow{d_0} \mathcal{N}(0, 1)$$

Для $H_1: p > q$ критерий $S = \{Z(x, y) > z_{1-\alpha}\}$

Дов. интервал для $p_1 - p_2$ равен $C = (\hat{p}_1 - \hat{p}_2 - z_{1-\alpha}\hat{\sigma}, 1)$,

где $\hat{\sigma} = \sqrt{\frac{f+g}{n^2} + \frac{(f-g)^2}{n^3}}$.

Пример

Два опроса с интервалом в полгода среди 1600 граждан Великобритании об одобрении премьер-министра.

	$Y_i = 1$	$Y_i = 0$	Σ
$X_i = 1$	$e = 794$	$f = 150$	944
$X_i = 0$	$g = 86$	$h = 720$	656
Σ	880	720	1600

H_0 : рейтинг не изменился

H_1 : рейтинг изменился $\implies pvalue = 2.8 \cdot 10^{-5}$

H_1 : рейтинг снизился $\implies pvalue = 1.4 \cdot 10^{-5}$

H_1 : рейтинг увеличился $\implies pvalue = 0.99999$

Доверительный интервал (0.0214, 0.0590)



Дов. интервал Уилсона для $p - q$

Независимые выборки

$$[\hat{p} - \hat{q} - \delta, \hat{p} - \hat{q} + \varepsilon],$$

$$\delta = z_{\alpha/2} \sqrt{\frac{\ell_1(1-\ell_1)}{n} + \frac{u_2(1-u_2)}{m}},$$

$$\varepsilon = z_{\alpha/2} \sqrt{\frac{u_1(1-u_1)}{n} + \frac{\ell_2(1-\ell_2)}{m}},$$

ℓ_1, u_1 — корни уравнения

$$|x - \hat{p}| = z_{\alpha/2} \sqrt{\frac{x(1-x)}{n}},$$

ℓ_2, u_2 — корни уравнения

$$|x - \hat{q}| = z_{\alpha/2} \sqrt{\frac{x(1-x)}{m}}.$$

Связные выборки

$$[\hat{p} - \hat{q} - \delta, \hat{p} - \hat{q} + \varepsilon],$$

$$\delta = \sqrt{d\ell_1^2 - 2\hat{\varphi}d\ell_1du_2 + du_2^2},$$

$$\varepsilon = \sqrt{du_1^2 - 2\hat{\varphi}du_1d\ell_2 + d\ell_2^2},$$

$$\hat{\varphi} = \begin{cases} \frac{eh - fg}{(e+f)(g+h)(e+h)(f+g)}; \\ 0, & \text{если знаменатель} = 0; \end{cases}$$

$$d\ell_1 = \hat{p} - \ell_1, \quad d\ell_2 = \hat{q} - \ell_2,$$

$$du_1 = u_1 - \hat{p}, \quad du_2 = u_2 - \hat{q},$$

ℓ_1, u_1 — корни уравнения

$$|x - \hat{p}| = z_{\alpha/2} \sqrt{\frac{x(1-x)}{n}},$$

ℓ_2, u_2 — корни уравнения

$$|x - \hat{q}| = z_{\alpha/2} \sqrt{\frac{x(1-x)}{n}}.$$



Нормальные выборки



Связные выборки

$$X_1, \dots, X_n \sim \mathcal{N}(a_1, \sigma_1^2)$$

$$Y_1, \dots, Y_n \sim \mathcal{N}(a_2, \sigma_2^2).$$

$$H_0: a_1 = a_2 \text{ vs. } H_1: a_1 \{<, \neq, >\} a_2$$

Сведение к задаче с одной выборкой:

Рассмотрим выборку $\delta_1, \dots, \delta_n$, где $\delta_i = X_i - Y_i$.

Тогда $H_0: E\delta_i = 0$ vs. $H_1: E\delta_i \{<, \neq, >\} 0$

Применяем критерий Вальда:

$$Z(X, Y) = \sqrt{n} \bar{\delta} / S_{\delta} \xrightarrow{d_0} \mathcal{N}(0, 1)$$

Почему не точный?

Если $X_i \sim \mathcal{N}(a_1, \sigma_1^2)$ и $Y_i \sim \mathcal{N}(a_2, \sigma_2^2)$ зависимы, то разность не обязана быть нормальной.



Связные выборки: обобщение

$$X_1, \dots, X_n$$

$$Y_1, \dots, Y_n.$$

$$H_0: EX_1 = EY_1 \text{ vs. } H_1: EX_1 \{<, \neq, >\} EY_1$$

Сведение к задаче с одной выборкой:

Рассмотрим выборку $\delta_1, \dots, \delta_n$, где $\delta_i = X_i - Y_i$.

Требование: $\delta_1, \dots, \delta_n$ — выборка с конечной дисперсией.

Тогда $H_0: E\delta_i = 0$ vs. $H_1: \delta_i \{<, \neq, >\} 0$

Применяем критерий Вальда:

$$Z(X, Y) = \sqrt{n} \bar{\delta} / S_{\delta} \xrightarrow{d_0} \mathcal{N}(0, 1)$$



Независимые выборки

$$X_1, \dots, X_n \sim \mathcal{N}(a_1, \sigma_1^2)$$

$$Y_1, \dots, Y_m \sim \mathcal{N}(a_2, \sigma_2^2).$$

Виды гипотез:

1. Равенство средних

$$H_0: a_1 = a_2 \quad \text{vs.} \quad H_1: a_1 \{<, \neq, >\} a_2$$

Способ зависит от доступной информации о дисперсиях.

2. Равенство дисперсий

$$H_0: \sigma_1 = \sigma_2 \quad \text{vs.} \quad H_1: \sigma_1 \{<, \neq, >\} \sigma_2$$

3. Однородность

$$H_0: (a_1, \sigma_1^2) = (a_2, \sigma_2^2)$$



Независимые выборки: среднее

$$X_1, \dots, X_n \sim \mathcal{N}(a_1, \sigma_1^2)$$

$$Y_1, \dots, Y_m \sim \mathcal{N}(a_2, \sigma_2^2).$$

$$H_0: a_1 = a_2$$

$$H_1: a_1 \{<, \neq, >\} a_2$$

Рассуждения:

$$\bar{X} \sim \mathcal{N}(a_1, \sigma_1^2/n)$$

$$\bar{Y} \sim \mathcal{N}(a_2, \sigma_2^2/m)$$

$$\bar{X} - \bar{Y} \overset{H_0}{\sim} \mathcal{N}(0, \sigma_1^2/n + \sigma_2^2/m)$$

Случай 1. σ_1 и σ_2 известны

Статистика критерия

$$Z(X, Y) = \frac{\bar{X} - \bar{Y}}{\sqrt{\sigma_1^2/n + \sigma_2^2/m}} \overset{H_0}{\sim} \mathcal{N}(0, 1)$$

Случай 2. $\sigma_1 = \sigma_2 = \sigma$ неизвестны

S_X^2, S_Y^2 — несмещ. оценки дисп.

Несмещенная оценка σ :

как взвешенное усреднение дисперсий

$$S_{tot}^2 = \frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}$$

Статистика критерия

$$T(X, Y) = \frac{\bar{X} - \bar{Y}}{S_{tot} \sqrt{1/n + 1/m}} \overset{H_0}{\sim} T_{n+m-2}$$

Случай 3. $\sigma_1 \neq \sigma_2$ и неизвестны

S_X^2, S_Y^2 — несмещ. оценки дисп.

$$T(X, Y) = \frac{\bar{X} - \bar{Y}}{\sqrt{S_X^2/n + S_Y^2/m}} \overset{H_0}{\underset{\text{прибл.}}{\sim}} T_v$$

$$v = \left(\frac{S_X^2}{n} + \frac{S_Y^2}{m} \right)^2 \left/ \left(\frac{S_X^4}{n^2(n-1)} + \frac{S_Y^4}{m^2(m-1)} \right) \right.$$



Независимые выборки: дисперсия

$$X_1, \dots, X_n \sim \mathcal{N}(a_1, \sigma_1^2)$$

$$Y_1, \dots, Y_m \sim \mathcal{N}(a_2, \sigma_2^2).$$

$$H_0: \sigma_1 = \sigma_2 \quad \text{vs.} \quad H_1: \sigma_1 \{<, \neq, >\} \sigma_2$$

Критерий Фишера

S_X^2, S_Y^2 — несмещенные оценки дисперсии

$$\frac{(n-1)S_X^2}{\sigma^2} \sim \chi_{n-1}^2, \quad \frac{(m-1)S_Y^2}{\sigma^2} \sim \chi_{m-1}^2$$

$$F(X, Y) = S_X^2 / S_Y^2 \stackrel{H_0}{\sim} F_{n-1, m-1}$$

Не устойчив к отклонениям от нормальности даже асимптотически, нужна строгая проверка нормальности критерием Шапиро-Уилка.



Независимые выборки: однородность

$$X_1, \dots, X_n \sim \mathcal{N}(a_1, \sigma_1^2)$$

$$Y_1, \dots, Y_m \sim \mathcal{N}(a_2, \sigma_2^2)$$

$$H_0: (a_1, \sigma_1^2) = (a_2, \sigma_2^2)$$

Идея:

1. Проверим гипотезу $H'_0: \sigma_1^2 = \sigma_2^2$
 \implies *двусторонний критерий Фишера*;
2. Если не отвергается, то проверим $H''_0: a_1 = a_2$
 \implies *двусторонний критерий Стьюдента (случай $\sigma_1 = \sigma_2$)*.

Как быть с уровнем значимости?

Применяем МПГ по методу Холма.



ВСЁ!