

МСПС | Дисперсионный анализ

(1)

Задачи дисперсионного анализа (однофакторности, в основном) чаще всего возникают в производстве, когда нужно сравнить выработанные изделия, к примеру, с контрольным образцом. Т.е. есть 2 выборки: одна - у рабочей совокупности, другая - контрольная. Ну и хотим, например, понять, есть ли эффект обработки, т.е. отличается ли рабочая выборка от контрольной по распределению.

Итак, пусть есть $X = (X_1, \dots, X_m)$ - м.о.р. с.в. с ф.р. $F(x)$; $Y = (Y_1, \dots, Y_n)$ - н.о.р. с.в. с ф.р. $G(x)$. Практически будем считать, что $F(x)$ и $G(x)$ - непрерывны.

Важным для практики случаем является случай нормально распределенных выборок, т.е. $X_i \sim N(\mu_1, \sigma_1^2)$, $Y_j \sim N(\mu_2, \sigma_2^2)$. Проверим эти выборки на однородность при условии, что эти выборки независимы.

Понятно, что $S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ несмещенно оцен. σ_1^2 и S_2^2

$$\frac{(n-1)S_1^2}{\sigma_1^2} \sim \chi_{n-1}^2, \text{ а } ES_2^2 = \sigma_2^2 \text{ и } \frac{(m-1)S_2^2}{\sigma_2^2} \sim \chi_{m-1}^2.$$

Математическое ожидание. Сл.в. ξ имеет распр. Фишера с k_1 и k_2 степенями свободы, если $\xi \stackrel{d}{=} \frac{\frac{1}{k_1} \cdot \zeta}{\frac{1}{k_2} \eta}$, где $\zeta \sim \chi_{k_1}^2$, $\eta \sim \chi_{k_2}^2$ и $\zeta \perp \eta$.

Обозн.: $\xi \sim F_{k_1, k_2}$.

Критерий Фишера. Если верна H' : $\sigma_1^2 = \sigma_2^2$, то $\frac{S_1^2}{S_2^2} \sim F_{n-1, m-1}$, тогда критерий: $\left\{ \frac{S_1^2}{S_2^2} > f_{1-\alpha} \right\}$, $f_{1-\alpha}$ - квантиль ур. $1-\alpha$ у $F_{n-1, m-1}$.

Здесь большая делится на меньшую.

Если крит. Фишера не отвергает H' , то проверим H'' : $\mu_1 = \mu_2$.

Поскольку $S_1^2 \perp S_2^2$, то $\frac{1}{2} [(n-1)S_1^2 + (m-1)S_2^2] \sim \chi_{n+m-2}^2$, кроме того, $E\zeta = k$, если $\zeta \sim \chi_k^2$. Тогда $S_{tot}^2 = \frac{(n-1)S_1^2 + (m-1)S_2^2}{n+m-2}$ несмещенно оценивает σ^2 . Далее, $\bar{X} - \bar{Y} \sim N(0, \sigma^2(\frac{1}{n} + \frac{1}{m}))$ и не зависит от S_{tot}^2 (т. об ортогональном разложении гаусс. вектора).

Математическое ожидание. Пусть z_0, \dots, z_n - неуп. $N(0, 1)$, тогда $\xi = \frac{z_0^2}{z_1^2 + \dots + z_n^2}$ имеет распр. Стьюдента с n степенями свободы. (обозн. $\xi \sim t_n$)

$$\text{Имеем } T = \frac{\bar{X} - \bar{Y}}{\sqrt{S_{tot}^2(\frac{1}{n} + \frac{1}{m})}} = \sqrt{\frac{nm}{n+m}} \frac{\bar{X} - \bar{Y}}{S_{tot}} \sim t_{n+m-2}$$

Это и есть критерий Стьюдента проверки однородности нормальных выборок.

Замечание 1) Оптимальен для норм. выборок, но его мощность сильно снижается при отклон. от этой нормальности.
2) Лучше брать \approx равные выборки.

Рассмотрим ещё несколько методов проверки однородности независимых выборок, если уже они не являются нормальными.

Критерий Смирнова

Определим $D_{n,m} = \sup_x |F_n^*(x) - G_m^*(x)|$, где $F_n^*(x)$ и $G_m^*(x)$ — эмпирические ф.р. выборок $\{X$ и Y соотв.

Теорема (Смирнов) Если верна $H_0: F=G$, то $\sqrt{\frac{nm}{n+m}} D_{n,m} \xrightarrow{d} K$, где K — сл.в., имеющая ф.р. Колмогорова, $F_K(z) = \sum_{i \in \mathbb{Z}} (-1)^i e^{-2i^2 z^2}$.
 Примечание: Хорошо работает при $n, m \geq 20$, при $n, m < 20$ искажает таблицу критич. знач. в обеих частях или с пом. моделирования.

Омега-квадрат критерий Рендла

Статистика критерия $\omega_{n,m}^2 = \int_{-\infty}^{+\infty} [F_n^*(x) - G_m^*(x)]^2 dH_{n+m}^*(x)$, где $H_{n+m}^*(x) = \frac{n}{n+m} F_n^*(x) + \frac{m}{n+m} G_m^*(x)$ — эмпир. ф.р. выборки $(X_1, \dots, X_n, Y_1, \dots, Y_m)$ (есть ещё один вид $\omega_{n,m}^2 = \frac{1}{n+m} \left[\frac{1}{6} + \frac{1}{m} \left\{ \sum_{i=1}^n (R_i - i)^2 \right\} + \frac{1}{n} \left\{ \sum_{j=1}^m (S_j - j)^2 \right\} \right] - \frac{2}{3}$, где R_i — ранг $X_{(i)}$, S_j — ранг $Y_{(j)}$ в обьед. вариационн. ряду).

Теорема (Рендла) Если $H_0: F=G$ верна, то при $\frac{n}{n+m} \rightarrow \lambda \in (0,1)$ $\frac{nm}{n+m} \omega_{n,m}^2 \xrightarrow{d} \chi^2$, где χ^2 имеет табличное распредел.

Замечание При малых n, m статистику $Z = \frac{nm}{n+m} \omega_{n,m}^2$ следует заменить на $Z^* = \frac{Z - EZ}{\sqrt{450Z}} + \frac{1}{6}$, где $EZ = \frac{1}{6} \left(1 + \frac{1}{n+m} \right)$, $DZ = \frac{1}{45} \left(1 + \frac{1}{n+m} \right) \left(1 + \frac{1}{n+m} - \frac{3}{4} \left(\frac{1}{n} + \frac{1}{m} \right) \right)$, это обесп. хорошую точность уже при $n, m \geq 7$.

Критерий рангов всех сумм Уилкоксона (или крит. Манна-Уитни-Уилкоксона, англ. обобщ. MWU)

Проверяем тут гипотезу о сдвиге: предполагаем, что $G(y) = F(y - \theta)$ и $H_0: \theta = 0$ (альтернативы бывают $H_+: \theta > 0$, $H_-: \theta < 0$, $H_1: \theta \neq 0$, для них всех критерий явл. составительным).

Рассм. обьединённую совокупность $(X_1, \dots, X_n, Y_1, \dots, Y_m)$, перейдём к рангам сл.в. в этой совокупности: $R(X_i) = R_i$ и $S_j = R(Y_j)$. $N = n+m$. Определим $W_{n,n} = \sum_{j=1}^m S_j$ — статистика ранг. сумм Уилкоксона.

Теорема При верной H_0 $\frac{W_{n,n} - m \cdot \frac{n+m+1}{2}}{\sqrt{nm \frac{n+m+1}{12}}} \xrightarrow{d} N(0,1)$.

Замечание При $n, m \leq 50$ статистика $W^* = \frac{W_{n,n} - m \cdot \frac{n+m+1}{2}}{\sqrt{nm \frac{n+m+1}{12}}}$ плохо аппроксимируется нормал. законом, в этом случае применять $\tilde{W} = \frac{1}{2} W^* \cdot \left(1 + \sqrt{\frac{N+2}{N-1 - (W^*)^2}} \right)$, где в кал. критич. значение брать $z_{1-\alpha} = (X_{1-\alpha} + Y_{1-\alpha})/2$, $X_{1-\alpha}$ и $Y_{1-\alpha}$ — квантили ур. $1-\alpha$ у $N(0,1)$ и T_{N-2} соотв. Если $n, m \geq 10$, то применять таблич. значение.

МСПС | Дисперсионный анализ

(3)

Представляет интерес также ситуация, когда выборки (X_1, \dots, X_n) и (Y_1, \dots, Y_n) явл. парными, т.е. (Z_1, \dots, Z_n) , где $Z_i = (X_i, Y_i)^T$, явл. выборкой из векторов, зависимость между X_i и Y_i неизвестна.

Например, такая задача возникает, когда мы делаем измерения на неск. приборе и делаем замеры дважды, пытаюсь понять, изменилось ли по-прежнему.

Из-за неизвест. зависимости между X_i и Y_i приходится пользоваться ранговыми методами, при этом будем проверять гипотезу об отсут. связи, как и в крит. Уилкоксона.

Рассмотрим $V_i = Y_i - X_i = \theta + \varepsilon_i$, где θ - интерес. нас эффект воздействия. Будем считать, что $P(\varepsilon_i \leq 0) = P(\varepsilon_i \geq 0) = \frac{1}{2}$ и распред. ε_i непрерывно. Кроме того, ε_i незав. по построению.

Критерий знаков

Рассм. $S = \sum_{i=1}^n I(V_i > 0)$. Тогда при верной $H_0: \theta = 0$ $\frac{S - \frac{n}{2}}{\sqrt{\frac{n}{4}}} \xrightarrow{d} N(0, 1)$.

При $n < 15$ нормал. аппроксимация работает плохо, то можно вычислить квантили из сум. : $P_0(S \geq k) = 2^{-n} \sum_{i=k}^n C_n^i$.

Замечание Если всё же распред. ε не явл. непрерывным и сферы значений V_i есть нулевые, то их надо отбросить и уменьшить n до числа ненул. эл-тов.

Критерий знаков с рангами Уилкоксона

Теперь вообще считаем, что $\varepsilon_1, \dots, \varepsilon_n$ - н.о.р. с симметр. непрер. распред.

Статистика: $T = \sum_{i=1}^n R_i \cdot I(V_i > 0)$, где R_i - ранг $|V_i|$ в возрастающ. ряду выборки $(|V_1|, \dots, |V_n|)$.

Теорема: при верной $H_0: \theta = 0$ $T^* = \frac{T - \frac{n(n+1)}{4}}{\sqrt{\frac{n(n+1)(2n+1)}{24}}} \xrightarrow{d} N(0, 1)$

При $n \leq 15$ следует пользоваться таблич. критич. значениями, т.к. нормальное приближение работает плохо.

Замечание 1) Поправка: $\tilde{T} = \frac{1}{2} T^* \left(1 + \sqrt{\frac{n-1}{n - (T^*)^2}} \right)$, уточняет

аппроксимацию при малых n . Тогда крит. знач. $\tilde{a}_\alpha = \frac{X_{1-\alpha} + Y_{1-\alpha}}{2}$, где $X_{1-\alpha}$ и $Y_{1-\alpha}$ - $(1-\alpha)$ -квантили у $N(0, 1)$ и t_{n-1} соотв.

2) Если всё-таки распред. ε_i не непрерывно, то

если $V_i = 0$ для нек. i , то отбрас. их и уменьш. n .

А если сферы $|V_i|$ есть равные, то используем средние ранги,

а дисперсия меняется так: $\frac{1}{24} [n(n+1)(2n+1) - \frac{1}{2} \sum_{k=1}^g \ell_k (\ell_k^2 - 1)]$,

где g - число групп совпадений, а ℓ_k - кол-во эл-тов в k -той группе.

3) Чтобы применить критерий ранговых знаков, нужно убедиться хотя бы в симметричности ε_i ; делается это так: должны примерно совпадать $z_i = V_{(i)} + \hat{\mu}$ и $\eta_i = V_{(n+1-i)} - \hat{\mu}$, $i = 1, \dots, [n/2]$, $\hat{\mu}$ - выбороч. медиана.

Дисперсионный анализ | МСРС |

(4)

Мы вот для корр. распределений проверили только гипотезу ортогонал., а иногда хочется и гипотезу о равенстве средних проверить.

Итак, пусть $X_1, \dots, X_n \sim N(\mu_1, \sigma_1^2)$, $Y_1, \dots, Y_m \sim N(\mu_2, \sigma_2^2)$ и выборки независимы.

Тогда при верной H_0 $T(X, Y) = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{m}}} \sim \overset{\text{таблицы}}{t_k}$, где

$$k = \frac{\left(\frac{s_1^2}{n} + \frac{s_2^2}{m}\right)^2}{\frac{s_1^4}{n^2(n-1)} + \frac{s_2^4}{m^2(m-1)}}$$

Критерий Асимики-Уэлша.

(С этой задачей связана известная проблема Беренда-Финша о несуществовании наилучшего критерия в этой задаче)

Пусть теперь наблюдения парные (ещё говорит, что выборки связанные) и $n=m$.

При верной $H_0: \mu_1 = \mu_2$ $T(X, Y) = \frac{\bar{X} - \bar{Y}}{S/\sqrt{n}} \sim t_{n-1}$, где

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (D_i - \bar{D})^2}, \quad D_i = X_i - Y_i.$$

(Это называется t -критерий Стьюдента для связанных выборок)

Кроме того, есть ~~несколько~~ множество непараметрических критериев, которые проверяют разные случаи гипотезы ортогональности — например, ~~и~~ гипотезу о равенстве средних (Манна-Уитни критерий) или гипотезу о равенстве дисперсий (критерий Зигеле-Тьюки).

Приведём критерий Зигеле-Тьюки.

Пусть выборки X и Y независимы, проверим гипотезу

$H_0: DX_1 = DY_1$ vs. $H_1: DX_1 \neq DY_1$.

Пусть $(Z_{(1)} \dots Z_{(n+m)})$ — вариационный ряд объединённой выборки $(X_1, \dots, X_n, Y_1, \dots, Y_m)$.

Статистика критерия $R_1(X, Y) = \sum_{i=1}^n \tilde{\text{rank}}(X_i)$, где $\tilde{\text{rank}}$ присваивается так: $\tilde{\text{rank}}(X_i) = \text{rank}(X_i)$ т.е. по Y_j ранги не суммируются.

$Z_{(1)}: Z_{(1)} \leq Z_{(2)} \leq Z_{(3)} \leq \dots \leq Z_{(n+m-2)} \leq Z_{(n+m-1)} \leq Z_{(n+m)}$

$\tilde{\text{rank}} \quad 1 \quad 4 \quad 5 \quad \dots \quad 6 \quad 3 \quad 2$

$R_1(X, Y)$ при верной H_0 имеет табличное распределение.