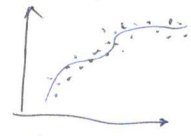
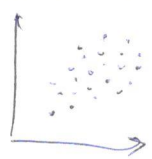


КПС | Корреляционный анализ

7

Пусть мы изучаем некоторые объекты, кот. обладают несколькими признаками. Часто возникает задача изучения взаимосвязей этих признаков (и дальнейшие задачи отнесения как более "высших" признаков и т.д.) Чаще всего встречаются 2 вида взаимосвязей (естественно, случаи упрощенные)

- а) объекты образуют "облако" точек т.е.
- б) объекты располож. в окрест. некоторой кривой



В случае а) оба признака являются количественными, и изучению подлежат уровни зависимости (корреляции) между ними. Случай б) соответствует "функционал." зависимости между признаками, именованной кривой. ~~В~~ этот случай изучается методами регресс. анализа. А сейчас будем изучать первый случай с помощью методов коррел. анализа

Сначала рассмотрим гипотезу независимости признаков z и η :
 $H_0: F_{z,\eta}(x,y) = F_z(x) \cdot F_\eta(y)$, где $F_{z,\eta}(x,y) - \text{Ф.р.}(z,\eta)$; $F_z(x)$ и $F_\eta(y) - \text{Ф.р.}$

Будем считать, что есть данные, и признаки z и η реализованы в виде выборок $(X_1 \dots X_n)$ и $(Y_1 \dots Y_n)$

Первая мысль - воспольз. ~~для~~ ^{для} ~~обычных~~ ^{коэфф. корреляции} $r = \frac{\text{cov}(X,Y)}{\sqrt{DX} \sqrt{DY}}$ (точно, его выбороч. характеристика). (Увы, если $r=0$, это не значит, что с.в. независимы. То же самое верно и для выборок. коэфф. корреляции. Подробнее об этом - на семинарах)

Коэфф. корреляции Пирсона

Опр. $\hat{r} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(\sum_{i=1}^n (X_i - \bar{X})^2 \cdot \sum_{i=1}^n (Y_i - \bar{Y})^2)^{1/2}}$

Св-ва 1) Если $EX_i^2 < +\infty$ и $EY_i^2 < +\infty$, то $\hat{r} \xrightarrow{P} r$

2) Если H_0 верна (а, точнее, H'_0 - гипотеза о некоррелируемости выборок), то $\frac{\hat{r} \sqrt{n-2}}{\sqrt{1-\hat{r}^2}} \rightarrow T_{n-2}$ (распред. Стьюдента)

- 3) Плохо реализуется на выбросы.
- 4) Для гаусс. выборок гипотезы H_0 и H'_0 эквивалентны, для них \hat{r} также верно и применяется.

Коэфф. корр. Спирмена

Пусть R_i - ранг наблюдения X_i в выборке $(X_1 \dots X_n)$, т.е. $X_{(R_i)} = X_i$, S_i - ранг Y_i .

Опр. Коэфф. корр. Спирмена: $\hat{r}_s = \frac{\sum_{i=1}^n (R_i - \bar{R})(S_i - \bar{S})}{\sqrt{\sum_{i=1}^n (R_i - \bar{R})^2 \cdot \sum_{i=1}^n (S_i - \bar{S})^2}}$, где $\bar{R} = \frac{1}{n} \sum_{i=1}^n R_i = \frac{n+1}{2} = \bar{S}$.

Св-ва ρ_S :

- 1) $\rho_S = 1 - \frac{6}{n^3 - n} \sum_{i=1}^n (R_i - S_i)^2$
- 2) При верной H_0 $E\rho_S = 0$; $D\rho_S = \frac{1}{n-1}$
- 3) $-1 \leq \rho_S \leq 1$ (из Коши-Бун.), причем $\rho_S = 1$, если $R_i = S_i$ и $\rho_S = -1$, если $R_i + S_i = n+1$
- 4) Распред. ρ_S не зависит от F_X и F_Y , его распред. табулировано.
- 5) При верной H_0 : $\frac{\rho_S}{\sqrt{D\rho_S}} \xrightarrow{d} N(0,1)$, приближением можно пользоваться при $n \geq 50$

6) Если $n \leq 50$, то лучше брать $\tilde{\rho}_S = \frac{1}{2} \rho_S (\sqrt{n-1} + \sqrt{\frac{n-2}{1-\rho_S^2}})$
 (поправка Шмидта). Тогда критерий: отвергнуть H_0 , если $\tilde{\rho}_S \notin (z_{\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}})$,
 где $z_\gamma = \frac{1}{2}(x_\gamma + y_\gamma)$, x_γ и y_γ - квантили $N(0,1)$ и T_{n-2} соответственно.
 7) Устойчив к выбросам (как и τ Кэжулла).

Коэфф. корреляции Кэжулла (τ Кэжулла)

Опр. Пары (X_i, Y_i) и (X_j, Y_j) - согласованы, если $\text{sign}(X_i - X_j)(Y_i - Y_j) = 1$.

Обозн. S - число согласов. пар, R - число несоглас., $T = S - R = \sum_{i < j} \text{sign}(X_i - X_j)(Y_i - Y_j)$

Легко видеть, что $-\frac{n(n-1)}{2} \leq T \leq \frac{n(n-1)}{2}$.

Опр. τ Кэжулла: $\tau = \frac{T}{\frac{n(n-1)}{2}}$.

Св-ва τ :

- 1) $\tau = 1 - \frac{4}{n(n-1)} R$
- 2) Если H_0 верна, то $E\tau = 0$ и $\frac{\tau}{\sqrt{D\tau}} \xrightarrow{d} N(0,1)$, где

$$D\tau = \frac{2(2n+5)}{9n(n-1)}$$

3) Верный ϕ -м: $\tau = 1 - \frac{4}{n^2 - n} \sum_{i < j} I(T_i > T_j)$ и $\rho_S = 1 - \frac{12}{n^3 - n} \sum_{i < j} (j-i) \cdot I(T_j < T_i)$

где T_i вводятся так: если $R_k = i$, то обознач. $T_i = S_k$.

4) ρ_S и τ сильно корр., $\rho(\rho_S, \tau) > 0,99$ даже при $n > 5$.

Обобщенный коэфф. корреляции

Удобен для реализации на комп. коэфф. корреляции.

Определенное отображение $C(X) = \|c_{ij}(X_i, X_j)\|_{i,j=1 \dots n}$, где $c_{ij} = -c_{ji}$ и $c_{ii} = 0$.

Опр. Обобщ. коэфф. корр.: $\hat{z} = \frac{\sum_{i < j} c_{ij}(X) \cdot c_{ij}(Y)}{\sqrt{\sum_{i < j} c_{ij}^2(X) \cdot \sum_{i < j} c_{ij}^2(Y)}}$

Тогда

- 1) при $c_{ij}(X) = X_j - X_i$ $\hat{z} = \rho$
- 2) при $c_{ij}(X) = R_j - R_i$ и $c_{ij}(Y) = S_j - S_i$ $\hat{z} = \rho_S$
- 3) при $c_{ij}(X) = \text{sign}(R_j - R_i)$ $\hat{z} = \tau$.

МОНС | Таблицы сопряженности |

(1)

Решаем задачу выв. статист. свдвд для сгруп. данных (т.е. разбитых на категории). Данные: $\|d_{ij}\|_{i=1..n, j=1..m}$, где d_{ij} - числовое значение об объектах с индек. i, j .

Пример

	мал	ср	больш
м	0.2	0.5	0.3
ж	1.6	2.5	0.4

 ← Потребление бюджет в день

Опишем 2 схемы, кот. приводят к таблицам сопряженности:

Схема I Пусть (d_{i1}, \dots, d_{im}) ; $i=1..n$ - независимые выборки из полиномиальных распр. с вероятн. q_{ij} и заданным числом набл. n_i , $\sum_{j=1}^m q_{ij} = 1$ и $\sum_{j=1}^m d_{ij} = n_i$. Т.е. наблюдаем n_i ^{незав.} экспериментов $P(\text{эксперимент завершится } j\text{-тым исходом}) = q_{ij}$; $d_{ij} = \# \{ \text{экспер., завершив. } j\text{-тым исходом} \}$ в i -той серии (Ср. с критерием хи-квадрат, то же самое делаем).

Эта схема позволяет проверить однородность нескольких орном. распр. (т.е. проверить гипотезу о их равенстве).

~~Проверяется~~ ставится эта гипотеза так: $H_I: q_{ij} = q_{.j} \forall i$, где $q_{.j} = \frac{1}{n} \sum_{i=1}^n q_{ij}$

Схема II Предполагается, что (d_{i1}, \dots, d_{im}) имеют полином. распр. с вероятн. (p_{i1}, \dots, p_{im}) и фиксир. числом наблюдений $N = \sum_{ij} d_{ij}$, $\sum_{ij} p_{ij} = 1$

Проверяем гипотезу независимости $H_{II}: p_{ij} = z_i \cdot s_j$, где $z_i = \sum_{j=1}^m p_{ik}$, $s_j = \sum_{i=1}^n p_{ij}$.

Почему эта схема для нас важна. Пусть есть 2 выборки $X(X_1, \dots, X_N)$ и $Y(Y_1, \dots, Y_N)$. Хотим проверить их независимость, т.е. что $F_{xy}(s, t) = F_X(s) \cdot F_Y(t)$.

~~Разделим R на 2 части~~ Пусть есть два разбиения R : $\{B_i\}_{i=1}^n$ и $\{C_j\}_{j=1}^m$. Обозначим $d_{ij} = \# \{ k: X_k \in B_i \text{ и } Y_k \in C_j \}$,

$p_{ij} := P(X_1 \in B_i, Y_1 \in C_j)$. Если выборки X и Y независимы, то $p_{ij} = P(X_1 \in B_i) \cdot P(Y_1 \in C_j) = (\sum_{k=1}^m p_{ik}) \cdot (\sum_{e=1}^n p_{ej}) = z_i \cdot s_j$!

Для проверки гипотез H_I и H_{II} применяется вариант крит. хи-квадрат: определим $\chi^2 = N \cdot \sum_{i=1}^n \sum_{j=1}^m \frac{(d_{ij} - \frac{n_i \cdot m_j}{N})^2}{n_i \cdot m_j}$, где $n_i = \sum_{k=1}^m d_{ik}$, $m_j = \sum_{e=1}^n d_{ej}$

Теорема $\chi^2 \xrightarrow{d} \chi_{(n-1)(m-1)}^2$ при $N \rightarrow +\infty$.

Замечание. С помощью критерия, построенного по этой теореме, действительно можно проверить гипотезу о независимости 2 выборок, но в отличие от критериев, построенных с помощью выворотных коэффициентов корреляции. Но проблема в том, что в теореме для хорошей работы критерия требуется, как и ранее, $\frac{n_i \cdot m_j}{N} \geq 5, \forall i, j$, а отвергнуть гипотезу о независимости с помощью коэф. корреляции можно на гораздо меньшем числе наблюдений.

МСПС1 Многомерная корреляция

Углубление стат. связи между $k \geq 3$ выборками.

Опр. (ранговой коэф. корреляции Кэндалла)

$$W = \frac{12}{k^2(n^3-n)} \sum_{i=1}^n \left(\sum_{j=1}^k R_{ij} - \frac{k(n+1)}{2} \right)^2, \text{ где } R_{ij} - \text{ранг (от 1 до } n) \text{ } i\text{-го}$$

эл-та в j -й выборке.

Св-ва:

1) $0 \leq W \leq 1$, причем $W=1 \Leftrightarrow$ все k ранжировок совпадают.

2) Обозначим через \bar{r}_S среднее арифм. коэфф. Симпсона по всем $\frac{k(k-1)}{2}$ парам выборок. Тогда

$$W = \frac{(k-1)\bar{r}_S + 1}{k}$$

3) $k(n-1)W \xrightarrow[n \rightarrow \infty]{d} J^2(n-1)$

Частная корреляция

Пример про труды и выс. през. брата.

Однако же, влияние неугнет. факторов на исследуемое переменное может искажать ~~эти~~ истинную связь между перемен., т.е. подсчеты могут приводить к ложным знач. парного коэфф. корр., а они могут быть независ. при угнетении этих факторов.

Опр. Частн. коэфф. корреляции между X и Y при искл. влияния сл. в. Z называется

$$r(X, Y | Z) = r_{XY|Z} = \frac{r(X, Y) - r(X, Z) \cdot r(Y, Z)}{\sqrt{(1 - r^2(X, Z))(1 - r^2(Y, Z))}}$$

К этой формуле приводит попытка исключить зависимость от Z , заменив X и Y сл. в. $X' = X - \alpha Z$, $Y' = Y - \beta Z$, коэф. корр. с Z , тогда "оставшаяся" корреляция ~~есть~~ есть корреляция Y' и X' .

Для попул. оценки $\hat{r}_{XY|Z}$ заменим все коэфф. корр. ~~на~~ на коэфф. корр. ~~выражения~~ Пирсона.

Св-во: Если X, Y, Z - выборки из незав. норм. законов, то $\hat{r}_{XY|Z}$ будет распредел. так же, как и \hat{r}_{XY} , но для выборок размера $n-1$. Тогда $\sqrt{n-1} \arcsin \hat{r}_{XY|Z} \xrightarrow{d} N(0, 1)$

Ранг. коэфф. корр. Кэндалла τ (в отл. от r_S) переносится на случай част. коэф. с пом. аналогич. ф-лы:

$$\tau_{XY|Z} = \frac{\tau_{XY} - \tau_{XZ} \tau_{YZ}}{\sqrt{(1 - \tau_{XZ}^2)(1 - \tau_{YZ}^2)}}.$$