

Лекция 9

Лука Yaroshevskiy

1 ноября

Содержание

1 Математическая регрессия	1
1.1 Метод наименьшие квадратов	1
2 Линейная парная регрессия	2
2.1 Геометрический смысл прямой линейной регрессии	2
2.2 Проверка гипотезы по значимости коэффициента линейной корреляции	3
2.3 Выборочное корреляционное отношение	3

1 Математическая регрессия

Определение. Регрессией X на Z называется функция $f(z) = E(X|Z = z)$.

$$x = f(z)$$

— уравнение регрессии, а ее график — линия регрессии

Пусть при n экспериментах, при значениях z_1, \dots, z_n случайной величины Z наблюдаемые значения x_1, \dots, x_n случайной величины X .

Обозначим ε разницу между наблюдаемыми и экспериментальными значениями X

$$\varepsilon_i = X_i - E(X|Z = z_i) = X_i - f(z_i)$$

Тогда $x_i = f(z_i) + \varepsilon_i$, $1 \leq i \leq n$, и ε_i — независимые случайные нормальные величины. $\varepsilon_i \in N(0, \sigma^2)$

$$a = 0, \text{ т.к. } E\varepsilon_i = E(X_i) - E(f(z_i)) = EX_i - E(X|Z = z) = 0$$

Цель: по данным значениям как можно более точно оценить функцию $f(z)$

Примечание. При этом предполагается (например из теории), что $f(z)$ — функция определенного типа, параметры которой не известны

1.1 Метод наименьшие квадратов

Метод наименьших квадратов состоит в выборе параметров функции $f(z)$, таким образом, чтобы минимизировать сумму квадратов ошибок

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (X_i - f(z_i))^2 \rightarrow \min$$

Определение. Пусть $\Theta = (\Theta_1, \dots, \Theta_n)$ — набор параметров функции $f(z)$.

$$\hat{\Theta} = \min \sum_{i=1}^n \varepsilon_i^2$$

2 Линейная парная регрессия

Пусть имеется *линейная регрессия*:

$$f(z) = a + bz$$

Тогда $X_i = a + bz_i + \varepsilon_i$, $1 \leq i \leq n$. Найдем оценки неизвестных параметров a и b методом наименьших квадратов (МНК).

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (X_i - a - bz_i)^2 \rightarrow \min$$

$$\frac{\partial}{\partial a} \sum_{i=1}^n \varepsilon_i^2 = 2 \sum_{i=1}^n (X_i - a - bz_i) \cdot (-1) = -2 \sum_{i=1}^n X_i + 2 \sum_{i=1}^n a + 2b \sum_{i=1}^n z_i = -2(n\bar{X} - na - bn\bar{z}) = 0$$

$$\frac{\partial}{\partial b} \sum_{i=1}^n \varepsilon_i^2 = 2 \sum_{i=1}^n (X_i - a - bz_i) \cdot (-z_i) = -2 \left(\sum_{i=1}^n X_i z_i - a \sum_{i=1}^n z_i - b \sum_{i=1}^n z_i^2 \right) = -2(n\bar{Xz} - an\bar{z} - bn\bar{z}^2) = 0$$

$$\begin{cases} n\bar{X} - na - bn\bar{z} = 0 \\ n\bar{Xz} - an\bar{z} - bn\bar{z}^2 \end{cases} \Leftrightarrow \begin{cases} \bar{X} - a - b\bar{z} = 0 \\ \bar{Xz} - z\bar{z} - b\bar{z}^2 \end{cases} \Leftrightarrow \begin{cases} a + b\bar{z} = \bar{X} \\ a\bar{z} + b\bar{z}^2 = \bar{Xz} \end{cases}$$

$$\begin{cases} a = \bar{X} - b\bar{z} \\ (\bar{X} - b\bar{z})\bar{z} + b\bar{z}^2 = \bar{Xz} \end{cases} \Leftrightarrow \begin{cases} a = \bar{x} - b\bar{z} \\ b(\bar{z}^2 - \bar{z}^2) = \bar{Xz} - \bar{X} \cdot \bar{z} \end{cases} \Leftrightarrow \begin{cases} a = \bar{X} - b\bar{z} \\ b = \frac{\bar{Xz} - \bar{X} \cdot \bar{z}}{\bar{z}^2 - \bar{z}^2} \end{cases}$$

$$\bar{X}_z = E(\bar{X}|Z = z) = f(z)$$

$$\bar{X}_z = a + bz$$

$$\bar{X}_z = \bar{X} - b\bar{z} + b\bar{z}$$

$$\bar{X}_z - \bar{X} = b(z - \bar{z})$$

$$\bar{X}_z - \bar{X} = \frac{\bar{Xz} - \bar{X} \cdot \bar{z}}{\hat{\sigma}_z^2} \cdot (z - \bar{z})$$

$$\bar{X}_z - \bar{X} = \frac{\bar{Zx} - \bar{X} \cdot \bar{z}}{\hat{\sigma}_z \cdot \hat{\sigma}_x} \cdot \frac{\hat{\sigma}_x}{\hat{\sigma}_z} \cdot (z - \bar{z})$$

$$\frac{\bar{X}_z - \bar{X}}{\hat{\sigma}_x} = \frac{\bar{Xz} - \bar{x} \cdot \bar{z}}{\hat{\sigma}_z \cdot \hat{\sigma}_x} \cdot \frac{z - \bar{z}}{\hat{\sigma}_z}$$

$$\frac{\bar{X}_z - \bar{X}}{\hat{\sigma}_x} = r_B \cdot \frac{\bar{z} - z}{\hat{\sigma}_z}$$

, где $r_B = \frac{\bar{Xz} - \bar{X} \cdot \bar{z}}{\hat{\sigma}_z \cdot \hat{\sigma}_x}$

2.1 Геометрический смысл прямой линейной регрессии

Прямая регрессии строится таким образом, чтобы сумма квадратов длин отрезков была наименьшей.

Определение. Выборочный коэффициент линейной корреляции

$$r_B = \frac{\bar{Xz} - \bar{X} \cdot \bar{z}}{\hat{\sigma}_x \cdot \hat{\sigma}_y}$$

Примечание. Отсюда видим, что выборочный коэффициент линейной корреляции является оценкой теоретического коэффициента полученного по методу моментов. Поэтому выборный коэффициент корреляции характеризует силу линейной связи между двумя случайными величинами. Знак r_B показывает является ли корреляция прямой или обратной

2.2 Проверка гипотезы по значимости коэффициента линейной корреляции

Пусть двумерная случайная величины распределена нормально. По выборке объема n вычислен r_B , а r — теоретический коэффициент линейной корреляции.

- $H_0 : r = 0$
- против $H_1 : r \neq 0$, т.е. коэффициент r_B статистически значимый

Теорема 2.1. Если H_0 верна, то статистика

$$K = \frac{r_B \sqrt{n-2}}{\sqrt{1-r_B^2}} \in T_{n-2}$$

— распределение Стьюдента с $n-2$ степенями свободы

Критерий: Пусть t_{α} — квантиль T_{n-2} уравнения значимости α

2.3 Выборочное корреляционное отношение

Пусть имеется k выборок случайной величины X , полученных при уровнях фактора z_1, \dots, z_k . Вычислены общая дисперсия D_O , D_M — межгрупповая дисперсия, D_B — внутригрупповая дисперсия. По теореме о разложении дисперсии $D_O = D_M + D_B$ Доделать