

Лекция 15. Элементы теории корреляции.

§1. Функциональная, статистическая и корреляционная зависимости.

Две случайные величины могут быть связаны функциональной зависимостью, т.е. изменение одной из них по определенному закону влечет изменение другой, или зависимостью другого рода, называется статистической, или быть независимыми.

Определение 1. Статистической называют зависимость, при которой изменение одной из величин влечет изменение распределения другой.

Определение 2. Если при изменении одной из величин изменяется среднее значение другой, то в этом случае статистическая зависимость.

Пример. называется корреляционной.

Пусть Y - урожай зерна, X - количество удобрений. С одинаковых по площади участков земли при равных количествах внесенных удобрений снимают различный урожай, т.е. Y не функция от X . Это объясняется случайными факторами (осадки, агротехника и т.д.). Как показывает опыт, средний урожай зависит от количество удобрений. Точка Y связана с X корреляционной зависимостью.

§2. Условные средние.

Пусть каждому X отвечает несколько Y . Например, значение X

$x_1 = 2$, $Y - y_1 = 5, y_2 = 6, y_3 = 10$. Найдем их среднее $\bar{y}_2 = \frac{5+6+10}{3} = 7$ - условное среднее.

Определение 1. Условными средним \bar{y}_x называется среднее арифметическое значение Y , соответствующее значению X .

Определение 2. Корреляционной зависимостью Y от X называют функциональную зависимость условной средней

$$\bar{y}_x \text{ от } x: \quad \bar{y}_x = f(x). \quad (1)$$

Уравнение (1) называется уравнением регрессии Y на X , функцию $f(x)$ называется регрессией Y на X , а ее график - линией регрессии Y на X .

Аналогично определяется корреляционная зависимость X от Y .

§3. Основные задачи теории корреляции.

Таких задач две.

Первая задача теории корреляции – установить форму корреляционной связи, т.е. вид функции регрессии (линейная, квадратичная, показательная и т.д.). Наиболее часто встречаются линейные. Если обе зависимости x от y и y от x линейны, то корреляция – линейная, в противном случае нелинейная.

Вторая задача теории корреляции – оценить тесноту (силу) корреляционной связи. Теснота корреляционной зависимости Y от X оценивается по величине

рассеяния значений Y вокруг словного среднего \bar{y}_x . Большое рассеяние говорит о слабой зависимости Y от X или об ее отсутствии. Малое рассеяние говорит о сильной связи, возможно и функциональной. Аналогично X от Y .

§4. Отыскание параметров выборочного уравнение прямой линии регрессии по не сгруппированным данным.

Допустим, корреляционная зависимость между Y , X линейная, тогда линии регрессии будут прямыми.

Для отыскания уравнений этих прямых проведено n независимых испытаний, в результате получено n пар чисел

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n). \quad (1)$$

Будем рассматривать их как случайную выборку из генеральной совокупности. Тогда величины и уравнения, найденные по этим данным, называются **выборочными**.

В простейшем случае, когда x соответствует $1/y$, можно искомое уравнение $\bar{y}_x = kx + b$ записать в виде $Y = kx + b$, где k – угловой коэффициент прямой линии регрессии и обозначать $k = \rho_{yx}$, тогда уравнение регрессии будет:
 $Y = \rho_{yx}x + b. \quad (2)$

Выберем параметры ρ_{yx} и b так, чтобы (1) были как можно ближе к (2) на плоскости X_0Y .

Назовем отклонением разность $Y_i - y_i$, ($i = \overline{1, n}$), где Y_i вычисляется по уравнению (2), y_i - наблюдаемая из (1). Подберем ρ_{yx} и b так, чтобы сумма квадратов отклонений была минимальной (т.е. это метод наименьших квадратов). Строим

функцию $F(\rho, b) = \sum_1^n (Y_i - y_i)^2$ или $F(\rho, b) = \sum_{i=1}^n (\rho_{yx}x_i + b - y_i)^2$. Для отыскания мини-

сумма составим $\frac{\partial F}{\partial \rho} = 2 \sum_1^n (\rho x_i + b - y_i)x_i = 0$ и отсюда получим систему 2-х линейных уравнений относительно ρ и b

$\left. \begin{aligned} (\sum x^2)\rho + (\sum x)b &= \sum xy \\ (\sum x)\rho + \sum nb &= \sum y \end{aligned} \right\}$, решив ее, найдем

$\rho = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}; b = \frac{\sum x^2 \sum y - \sum x \sum yx}{n \sum x^2 - (\sum x)^2}$ получим искомое уравнение

$y = \rho x + b.$

Пример. Найти выборочное уравнение прямой линии регрессии Y на X по данным наблюдений $n = 5$.

x	1	1,5	3	4,5	5
y	1,25	1,4	1,5	1,75	2,25

Решение: составим расчетную таблицу

	x_i	y_i	x_i^2	$x_i y_i$
1	1	1,25	1	1,25
2	1,5	1,4	2,25	2
3	3	1,5	9	4,5
4	4,5	1,75	20,25	4,875
5	5	2,25	25	11,25
	$\sum x_i = 15$	$\sum y_i = 8,15$	$\sum x_i^2 = 57,5$	$\sum x_i y_i = 26,975$

$$\rho_{yx} = \frac{5 \cdot 26,975 - 15 \cdot 8,15}{5 \cdot 57,5 - 15^2} = 0,202$$

По известным формулам вычислим

$$b = \frac{57,5 \cdot 8,15 - 15 \cdot 26,975}{62,5} = 1,024$$

Напиши искомое уравнение регрессии

$$Y = 0,202x + 1,024.$$

Для того, чтобы получить представление, насколько хорошо вычисленные по этому уравнению значения Y_i согласуются с наблюдаемыми значениями y_i , найдем отклонение $Y_i - y_i$.

x_i	Y_i	y_i	$Y_i - y_i$
1	1,226	1,25	-0,024
1,5	1,327	1,4	-0,073
3,00	1,630	1,5	0,130
4,5	1,933	1,75	0,083
5	2,034	2,25	-0,216

Из таблицы видно, что не все отклонения малы. Это объясняется числом наблюдений.

§5. Корреляционная таблица.

При большом числе наблюдений одно и то же значение x может встретиться n_x раз, одно и то же значение y n_y раз, одна и та же пара чисел (x, y) может наблюдаться n_{xy} раз. Поэтому данные наблюдений группируются, т.е. подсчитывают частоты n_x , n_y , n_{xy} . Все сгруппированные данные записывают в виде таблицы, которую называют корреляционной.

Пример.

X	10	20	30	40	n_x
Y					

0,4	5	-	7	14	26
0,6	-	2	6	4	12
0,8	3	19	-	-	22
n_x	8	21	13	18	$n = 60$

§6. Отыскание параметров выборочного уравнения прямой линии регрессии по сгруппированным данным. Выборочный коэффициент корреляции.

Для определения параметра уравнение прямой линии регрессии Y на X была получена система уравнений

$$\begin{cases} (\sum x^2)\rho_{yx} + (\sum x)b = \sum xy \\ (\sum x)\rho_{yx} + nb = \sum y \end{cases} \quad (*)$$

ρ_{yx} - коэффициент регрессии.

Предполагалось, что значения X и соответствующие им значения Y наблюдались по одному разу.

Теперь же допустим, что получено большое число данных (≈ 50 наблюдений), среди них есть повторяющиеся и они сгруппированы в виде корреляционной таблицы. Запишем систему (*) так, чтобы она отражала данные корреляционной таблицы

$$\sum x = n\bar{x} \quad (\text{следствие из } \bar{x} = \frac{\sum x^2}{n})$$

$$\sum y = n\bar{y}$$

$$\sum x^2 = n\bar{x}^2 \quad (\text{из } \bar{x}^2 = \frac{\sum x^2}{n})$$

$$\sum xy = \sum n_{xy}x_y \quad (\text{учтено, что } (x_y) \text{ наблюдается } n_{xy} \text{ раз}),$$

тогда (*) примет вид

$$\begin{cases} (n\bar{x}^2)\rho_{yx} + (n\bar{x})b = \sum n_{xy}xy \\ n(\bar{x})\rho_{yx} + b = \sum \bar{y} \end{cases} \quad (**)$$

Решив эту систему, найдем ρ_{yx} и b и получим $\bar{y}_x = \rho_{yx}x + b$.

Однако, удобнее, введя новую величину - коэффициент корреляции, написать уравнение регрессии в ином виде. Сделаем это $b = \bar{y} - \rho_{yx}\bar{x}$, подставив в уравнение регрессии, получим

$$\begin{aligned} \bar{y}_x &= \rho_{yx}x + \bar{y} - \rho_{yx}\bar{x} \\ \bar{y}_x - \bar{y} &= \rho_{yx}(x - \bar{x}) \end{aligned} \quad (***)$$

найдем из (**) ρ_{yx} , учитывая, что $\bar{x}^2 - (\bar{x})^2 = \sigma_x^2$

$$\rho_{yx} = \frac{\sum n_{xy}xy - n\bar{x}\bar{y}}{n(\bar{x}^2 - (\bar{x})^2)} = \frac{\sum n_{xy}xy - n\bar{x}\bar{y}}{n\sigma_x^2}$$

$$\rho_{yx} \cdot \frac{\sigma_x}{\sigma_y} = \frac{\sum n_{xy} xy - n\bar{x}\bar{y}}{n\sigma_x^2}$$

r_B – выборочный коэффициент корреляции

$$\rho_{yx} = r_B \cdot \frac{\sigma_y}{\sigma_x}$$

$\bar{y}_x - \bar{y} = r_B \frac{\sigma_y}{\sigma_x} (x - \bar{x})$ - уравнение прямой линии регрессии Y на X .

§7. Свойства выборочного коэффициента корреляции.

Цель введения этой величины – оценка тесноты линейной корреляционной зависимости. Это следует из его свойств.

$$S_y = D_y(1 - r_B^2), \quad S_x = D_x(1 - r_B^2),$$

где S_y - дисперсия наблюдавшихся значений y вокруг соответствующих условных средних \bar{y}_x ;

D_y - дисперсия наблюдений вокруг общей средней \bar{y} . Аналогично S_x и D_x .

1. абсолютная величина r_B не превосходит 1.

Доказательство: $S_y = D_y(1 - r_B^2) \geq 0$

$$1 - r_B^2 \geq 0, \quad -1 \leq r_B \leq 1 \text{ или } |r_B| \leq 1.$$

2. если $r_B = 0$ и выборочные линии регрессии – прямые, то X и Y не связаны линейной корреляционной зависимостью.

Доказательство: если $r_B = 0$

$$\bar{y}_x - \bar{y} = r_B \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$\bar{y}_x - \bar{y} = 0 \Rightarrow \bar{y}_x = \bar{y}$$

\bar{y}_x - сохраняет постоянное значение.

3. Если $|r_B| = 1$, то $S_y = D_y(1 - r_B^2) = 0$.

4. С возрастанием $|r_B|$ линейная корреляционная зависимость становится более тесной и при $|r_B| = 1$ переходит в функциональную зависимость.

Из этих всех свойств следует смысл r_B : выборочный коэффициент корреляции характеризует тесноту связи между количественными признаками в выборке, чем ближе $|r_B| \rightarrow 1$, тем связь сильнее, чем ближе $|r_B| \rightarrow 0$, тем слабее.

Для оценки коэффициента корреляции r_B нормально распределенной генеральной совокупности ($n \geq 50$) можно воспользоваться формулой

$$r_B - 3 \frac{1 - r_B}{\sqrt{4}} \leq r_2 \leq r_B + 3 \frac{1 + r_B}{\sqrt{4}}$$

8. Метод четырех полей вычисления коэффициента корреляции.

Дана корреляционная таблица. Вычислить r_B -?

Перейдем к угловым вариантам $u_i = \frac{x_i - a}{n_1}$, $v_j = \frac{y_j - a}{n_2}$, тогда $r_B = \frac{\sum u_n u v - n \bar{u} \bar{v}}{n \sigma_u \sigma_v}$.

\bar{u} , \bar{v} , σ_u и σ_v - вычисляются по методу произведений.

Для вычисления $\sum n_{uv} u \cdot v$ применяется метод 4-х полей. Название метода связано с тем, что строка и столбец, пересекающиеся в клетке, содержащей наибольшую частоту, делят корреляционную таблицу на H-части, которые называются полями. Рассмотрим метод на примере.

X	10	20	30	40	50	60	n_y
Y							
15	5	7	-	-	-	-	12
25	-	20	23	-	-	-	43
35	-	-	30	47	2	-	79
45	-	-	10	11	20	6	47
55	-	-	-	9	7	3	19
n_x	5	27	63	67	29	9	$n = 200$

В качестве C_1 взят вариант 40, имеющий наибольшую частоту

$$C_2 = 35, \quad h_1 = 20 - 10 = 10 \left. \vphantom{C_2} \right\} u = \frac{x - 40}{10}, \quad v = \frac{y - 35}{10}$$

u	-3	-2	-1	0	1	2	n_u	I	II
v									
-2	5 6	7 4					12	12	
-1		20 2	23 1				43	63	
0			30	47	2	-	79	III	IV
1			10 -1	11	20 11	6 2	47	-10	32
2				9	7 2	3 4	19	-	26
n_u	5	27	63	67	29	9	200		
I	30	68	23	II			121		
III			-10	IV					

u	-3	-2	-1	1	2	I	II
v							
-2	5 6	7 4				58	
-1		2	1			63	

		20	23				
1			-1	1	2	-10	32
			10	20	6		
2				2	4		26
				7	3		
I	30	68	23			121	
III			-10	34	24	-10	58

$$121+58=169$$

$\bar{u}, \bar{v}, \sigma_u$ и σ_v ,

$$\sigma_u = \sqrt{\bar{u}^2 - (\bar{u})^2}, \quad \sigma_v = \sqrt{\bar{v}^2 - (\bar{v})^2}$$

$$\bar{u} = \frac{\sum n_u u}{n} = \frac{5(-3) + 27(-2) + 63(-1) + 29 \cdot 1 + 9 \cdot 2}{200} = 0,425$$

$$\bar{v} = \frac{\sum n_v v}{n} = \frac{12(-1) + 43(-2) + 47 \cdot 1 + 19 \cdot 2}{200} = 0,09$$

$$\bar{u}^2 = \frac{5 \cdot 9 + 27 \cdot 4 + 63 \cdot 1 + 29 \cdot 1 + 9 \cdot 4}{200} = 1,405$$

$$\sigma_u = \sqrt{\bar{u}^2 - (\bar{u})^2} = \sqrt{1,405 - 0,425^2} = 1,106$$

$$\sigma_v = 1,209$$

$$\sum n_{uv} uv = 121 - 10 + 58 = 169.$$

$$r_B = \frac{\sum n_{uv} uv - n\bar{u}\bar{v}}{n\sigma_u\sigma_v} = \frac{169 - 200(-0,425) \cdot 0,09}{200 \cdot 1,106 \cdot 1,203} = 0,603$$

$$r_B = 0,603.$$

§9. Пример на отыскание выборочного уравнения прямой линии регрессии.

Так как имеются r_B и $\bar{u}, \bar{v}, \sigma_u$ и σ_v , то $\sigma_x = h_1\sigma_u$, $\sigma_y = h_2\sigma_v$, $\bar{x} = \bar{u}h_1 + C_1$, $\bar{y} = \bar{v}h_2 + C_2$.

Пример. Общий вид уравнения

$$\bar{y}_x - \bar{y} = r_B \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$\bar{x} = \bar{u}h_1 + C_1 = -0,425 \cdot 10 + 40 = 35,75$$

$$\bar{y} = \bar{v}h_2 + C_2 = 0,09 \cdot 10 + 35 = 35,9$$

$$\sigma_x = \sigma_u h_1 = 1,106 \cdot 10 = 11,06$$

$$\sigma_y = \sigma_v h_2 = 1,209 \cdot 10 = 12,09$$

$$\bar{y}_x - 35,9 = 0,603 \cdot \frac{12,09}{11,06} (x - 35,75)$$

$$\bar{y}_x = 0,659x + 12,34$$