

Поставим следующую задачу: произведено n взаимно независимых экспериментов с двумерной случайной величиной (X, Y) . Получены пары чисел $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$. По имеющимся данным определить ковариацию $\text{cov}(X, Y) = E((X - EX) \cdot (Y - EY)) = E(XY) - (EX) \cdot (EY)$.

Напомним в очередной раз, что это невозможно. И единственное, что удаётся сделать, это найти оценку, то есть формулу, определяющую значение $\text{cov}(X, Y)$ по числам $X_1, Y_1, X_2, Y_2, \dots, X_n, Y_n$ наилучшим образом с некоторой точки зрения. Самое простое решение заключается в том, что формулы берутся совпадающими с соответствующими им формулами для двумерной дискретной случайной величины, принимающей возможные значения $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ с равными вероятностями (равными $1/n$). Таблица распределения такой двумерной дискретной случайной величины имеет вид:

$x_i:$	X_1	X_2	\dots	X_n
$y_i:$	Y_1	Y_2	\dots	Y_n
$p_i:$	$1/n$	$1/n$	\dots	$1/n$

Оценки математических ожиданий, называемые «выборочные средние» известны:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{и} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

Для ковариации $\text{cov}(X, Y)$ хочется предложить оценку $\frac{1}{n} \sum_{i=1}^n ((X_i - EX)(Y_i - EY))$, заменив в ней математические ожидания EX и EY доступными для вычислений средними значениями $\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i$ и $\bar{y} = \frac{1}{n} \sum_{i=1}^n Y_i$. Тогда для ковариации $\text{cov}(X, Y)$ получится оценка

$$\frac{1}{n} \sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right) \left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right) \right).$$

Проверим эту оценку на выполнение свойства несмещённости: найдём

$$\begin{aligned} E \left(\frac{1}{n} \sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right) \left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right) \right) \right) &= \frac{1}{n} \cdot E \left(\sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right) \left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right) \right) \right) = \\ &= \frac{1}{n} \cdot \sum_{i=1}^n E \left(X_i Y_i - X_i \frac{1}{n} \left(\sum_{j=1}^n Y_j \right) - Y_i \frac{1}{n} \left(\sum_{j=1}^n X_j \right) + \frac{1}{n^2} \left(\sum_{j=1}^n X_j \right) \left(\sum_{j=1}^n Y_j \right) \right) = \end{aligned}$$

Это сумма n одинаковых слагаемых, поскольку все X_i имеют одинаковые распределения, такие же, как X_1 и, соответственно, все Y_i имеют одинаковые распределения, такие же, как Y_1 .

$$\begin{aligned} &= \frac{1}{n} \cdot n E \left(X_1 Y_1 - X_1 \frac{1}{n} \left(\sum_{j=1}^n Y_j \right) - Y_1 \frac{1}{n} \left(\sum_{j=1}^n X_j \right) + \frac{1}{n^2} \left(\sum_{j=1}^n X_j \right) \left(\sum_{j=1}^n Y_j \right) \right) = \\ &= E(X_1 Y_1) - E \left(X_1 \frac{1}{n} \left(\sum_{j=1}^n Y_j \right) \right) - E \left(Y_1 \frac{1}{n} \left(\sum_{j=1}^n X_j \right) \right) + E \left(\frac{1}{n^2} \left(\sum_{j=1}^n X_j \right) \left(\sum_{j=1}^n Y_j \right) \right) = \\ &= E(X_1 Y_1) - \frac{1}{n} \cdot E \left(X_1 \left(\sum_{j=1}^n Y_j \right) \right) - \frac{1}{n} \cdot E \left(Y_1 \left(\sum_{j=1}^n X_j \right) \right) + \frac{1}{n^2} \cdot E \left(\left(\sum_{j=1}^n X_j \right) \left(\sum_{j=1}^n Y_j \right) \right) = \end{aligned}$$

Теперь нужно понять следующее. Результаты разных испытаний независимые. А математическое ожидание независимых случайных величин равно произведению их математических ожиданий. Например, при $j > 1$ верно $E(X_1 Y_j) = EX_1 \cdot EY_j$. Поэтому для

продолжения выкладок нужно отделить слагаемые, где есть произведения результатов испытаний с разными индексами от произведений результатов с одинаковыми индексами.

$$= E(X_1 Y_1) - \frac{1}{n} \cdot E\left(X_1 \left(Y_1 + \sum_{j=2}^n Y_j\right)\right) - \frac{1}{n} \cdot E\left(Y_1 \left(X_1 + \sum_{j=2}^n X_j\right)\right) + \frac{1}{n^2} \cdot E\left(\sum_{j=1}^n \sum_{k=1}^n X_j Y_k\right) =$$

Суммы умножились обычным образом (умножается каждое слагаемое первой суммы на каждое слагаемое второй суммы).

$$= E(X_1 Y_1) - \frac{1}{n} \cdot E\left(X_1 Y_1 + \sum_{j=2}^n X_1 Y_j\right) - \frac{1}{n} \cdot E\left(X_1 Y_1 + \sum_{j=2}^n X_j Y_1\right) + \frac{1}{n^2} \cdot E\left(\sum_{j=1}^n X_j Y_j + \sum_{j=1}^n \sum_{k=1, k \neq j}^n X_j Y_k\right) =$$

В двойной сумме необходимо было отделить друг от друга слагаемые с одинаковыми индексами от слагаемых с разными индексами.

$$\begin{aligned} &= E(X_1 Y_1) - \frac{1}{n} \cdot \left(E(X_1 Y_1) + E\sum_{j=2}^n X_1 Y_j\right) - \frac{1}{n} \cdot \left(E(X_1 Y_1) + E\sum_{j=2}^n X_j Y_1\right) + \frac{1}{n^2} \cdot \left(E\sum_{j=1}^n X_j Y_j + E\sum_{j=1}^n \sum_{k=1, k \neq j}^n X_j Y_k\right) = \\ &= E(X_1 Y_1) - \frac{1}{n} \cdot \left(E(X_1 Y_1) + \sum_{j=2}^n E(X_1 Y_j)\right) - \frac{1}{n} \cdot \left(E(X_1 Y_1) + \sum_{j=2}^n E(X_j Y_1)\right) + \frac{1}{n^2} \cdot \left(\sum_{j=1}^n E(X_j Y_j) + \sum_{j=1}^n \sum_{k=1, k \neq j}^n E(X_j Y_k)\right) = \\ &= E(X_1 Y_1) - \frac{1}{n} \cdot \left(E(X_1 Y_1) + \sum_{j=2}^n E(X_1)E(Y_j)\right) - \frac{1}{n} \cdot \left(E(X_1 Y_1) + \sum_{j=2}^n E(X_j)E(Y_1)\right) + \\ &+ \frac{1}{n^2} \cdot \left(\sum_{j=1}^n E(X_j Y_j) + \sum_{j=1}^n \sum_{k=1, k \neq j}^n E(X_j)E(Y_k)\right) = \end{aligned}$$

В очередной раз обратим внимание, что в каждой из сумм все слагаемые одинаковые. В частности, $E(X_j) = E(X_1)$, и $E(Y_k) = E(Y_1)$.

Аккуратно подсчитаем количества слагаемых в каждой из сумм.

$$\begin{aligned} &= E(X_1 Y_1) - \frac{1}{n} \cdot (E(X_1 Y_1) + (n-1)E(X_1)E(Y_1)) - \frac{1}{n} \cdot (E(X_1 Y_1) + (n-1)E(X_1)E(Y_1)) + \\ &+ \frac{1}{n^2} \cdot (nE(X_1 Y_1) + (n^2 - n)E(X_1)E(Y_1)) = \\ &= E(X_1 Y_1) - \frac{2}{n} \cdot E(X_1 Y_1) - \frac{2(n-1)}{n} E(X_1)E(Y_1) + \frac{1}{n} \cdot E(X_1 Y_1) + \frac{n-1}{n} E(X_1)E(Y_1) = \\ &= E(X_1 Y_1) - \frac{1}{n} \cdot E(X_1 Y_1) - \frac{n-1}{n} E(X_1)E(Y_1) = \frac{n-1}{n} \cdot E(X_1 Y_1) - \frac{n-1}{n} E(X_1)E(Y_1) = \\ &= \frac{n-1}{n} \cdot (E(X_1 Y_1) - E(X_1)E(Y_1)) = \frac{n-1}{n} \cdot \text{cov}(X_1, Y_1) = \frac{n-1}{n} \cdot \text{cov}(X, Y) \neq \text{cov}(X, Y) \end{aligned}$$

Оценка не прошла тест на несмещённость. В то же время ясно, как исправить эту оценку, чтобы её математическое ожидание было равно $\text{cov}(X, Y)$. Нужно умножить прежнюю оценку на дробь $\frac{n}{n-1}$. Таким образом «выборочную ковариацию» разумнее вычислять по формуле

$$\text{cov}(X, Y) = \frac{n}{n-1} \cdot \frac{1}{n} \sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right) \left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right) \right) = \frac{1}{n-1} \sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right) \left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right) \right).$$

Такая оценка уже будет несмещённой.

Обычно расчётчики найдя ковариацию не останавливаются и вычисляют корреляцию

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{DX} \cdot \sqrt{DY}}.$$

Подставив найденные оценки для ковариации и дисперсий получим расчетную формулу для оценки корреляции

$$r_{XY} = \frac{\frac{1}{n-1} \sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right) \left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right) \right)}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right)^2 \right)} \cdot \sqrt{\frac{1}{n-1} \sum_{i=1}^n \left(\left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right)^2 \right)}}$$

Дробь можно сократить

$$r_{XY} = \frac{\sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right) \left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right) \right)}{\sqrt{\sum_{i=1}^n \left(\left(X_i - \frac{1}{n} \sum_{j=1}^n X_j \right)^2 \right)} \cdot \sqrt{\sum_{i=1}^n \left(\left(Y_i - \frac{1}{n} \sum_{j=1}^n Y_j \right)^2 \right)}}$$