



Корреляционнорегрессионный анализ данных. Парная регрессия

Кагирова Мария Вячеславовна, доцент, кандидат экономических наук, кафедра статистики и кибернетики РГАУ-МСХА им. К. А. Тимирязева





Вопросы:

- Показатели связи между признаками
- Парная линейная регрессия
- Оценка достоверности модели регрессии.
- Нелинейная регрессия

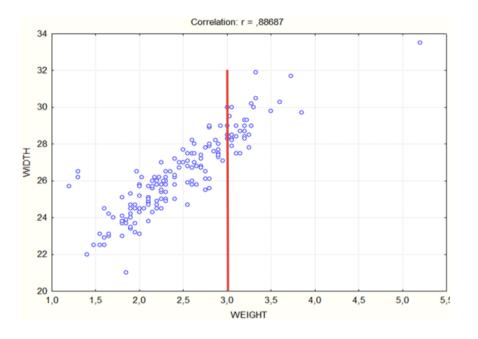






Функциональные связи являются «жесткими», четко определяющими значение одной переменной (признака) при изменении другой переменной (признака).

Статистические связи формируются под влиянием случайных факторов. При этом конкретному значению одной переменной может соответствовать не одно значение другой переменной, а множество, т.е. некоторое распределение второй переменной







Выборочная ковариация (Cov(x,y)) является мерой взаимосвязи между двумя переменными. Основой формирования данного показателя является оценка взаимной изменчивости переменных, выраженной в произведении отклонений каждой из них от среднего уровня.

Если Cov(x,y)>0, то связь прямая, если Cov(x,y)<0, то связь обратная, Cov(x,y)=0 — связь отсутствует.

$$Cov(x, y) = \frac{1}{n} \sum_{1}^{n} (x - \overline{x})(y - \overline{y})$$

$$Cov(x, y) = \frac{\sum xy}{n} - x \cdot y$$





Коэффициент корреляции
$$0<|r_{xy}|<1$$
: $0-0,3-$ связь слабая; $0,3-0,5-$ умеренная; $0,5-0,7-$ средняя; $0,7-0,99-$ тесная.

$$r_{xy} = \frac{Cov(x, y)}{\sqrt{Var(x)Var(y)}}$$



Показатели связи между признаками

	У	х	Z
у	1,000		
х	-0,470	1,000	
Z	0,055	0,825	1,000

Матрица парных коэффициентов корреляции

$$r_{xy.z} = \frac{r_{xy} - r_{xz}r_{yz}}{\sqrt{(1 - r_{xz}^2)(1 - r_{yz}^2)}} = \frac{-0.47 - 0.825 \cdot 0.055}{\sqrt{(1 - 0.825^2)(1 - 0.055^2)}} = -0.9146$$

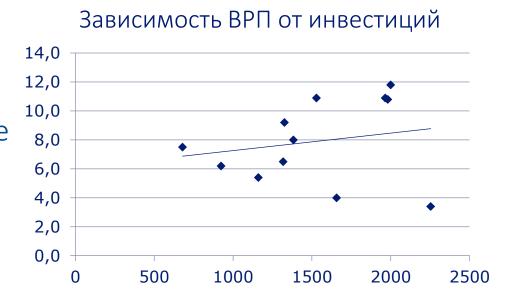




$$y = \alpha + \beta x + u$$

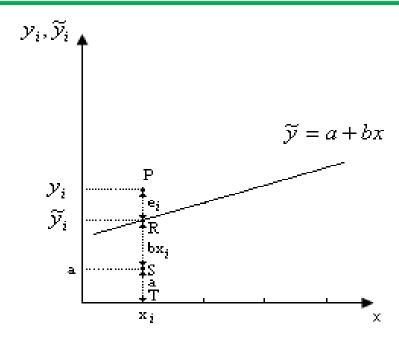
Величина у, рассматриваемая как зависимая переменная, состоит из двух составляющих:

- неслучайной составляющей $\alpha + \beta x$, где x выступает как объясняющая (или независимая) переменная, а постоянные величины α и β как параметры уравнения;
- случайного члена и.





Парная линейная регрессия



а – условное начало, характеризует значение зависимой переменной у при равенстве независимой переменной х нулю

b — коэффициент полной регрессии показывает, на сколько в среднем изменится значение зависимой переменной у (в своих единицах измерения) при изменении независимой переменной х на единицу (в своих единицах измерения)





Коэффициент детерминации

характеризует долю вариации зависимой переменной у, объясненную влиянием фактора х, включенного в уравнение регрессии.

Позволяет оценить качество модели регрессии. Чем ближе значение ${\sf R}^2$ к 1, тем выше качество модели.

$$R^{2} = \frac{Var(\tilde{y})}{Var(y)}$$

$$R^2 = 1 - \frac{Var(e)}{Var(y)}$$



Оценка достоверности модели регрессии

F-тест

Пригодность линии регрессии для прогноза зависит от того, какая часть общей вариации признака у приходится на объясненную вариацию. Очевидно, что если сумма квадратов отклонений, обусловленная регрессией, будет больше остаточной суммы квадратов, то уравнение регрессии статистически значимо и фактор *х* оказывает существенное воздействие на результат *у*





• Выдвигаются гипотезы:

$$H_0$$
: $\sigma^2_{perp} = \sigma^2_{oct}$
 H_a : $\sigma^2_{perp} \neq \sigma^2_{oct}$

- Выбирается уровень значимости критерия.
- Производится разложение общего объема вариации:
- Определяется число степеней свободы
- Рассчитываются выборочные несмещенные оценки дисперсий:
- Определяется фактическое значение F-критерия Фишера:
- Определяется критическое (табличное) значение критерия:
- Делается статистический вывод:

$$F$$
факт.≤ F табл. \Rightarrow H_0 ($\sigma^2_{perp} = \sigma^2_{oct}$)
 F факт.> F табл. \Rightarrow H_a ($\sigma^2_{perp} \neq \sigma^2_{oct}$)

$$S^{2}_{perp.} = \frac{W_{perp.}}{1}$$
$$S^{2}_{e} = \frac{W_{e}}{n-2}$$

$$F = \frac{S^2_{perp.}}{S^2_{e}}$$

$$F_{\alpha;v}$$
 perp. v oct.





t-Tect

• Формулируются рабочая и альтернативная гипотезы:

$$H_0: \alpha = 0; \beta = 0; \rho = 0$$
 $H_a: \alpha \neq 0; \beta \neq 0; \rho \neq 0$

- Выбирается уровень значимости критерия.
- Рассчитываются средние ошибки выборочных характеристик:

$$m_{a} = \sqrt{S_{e}^{2} \frac{\sum X_{i}^{2}}{n \sum X_{i} - \bar{X}^{2}}} = \sqrt{S_{e}^{2} \frac{\bar{X}^{2}}{n \sigma_{x}^{2}}} \qquad m_{b} = \sqrt{\frac{S_{e}^{2}}{\sum X_{i} - \bar{X}^{2}}} = \sqrt{\frac{S_{e}^{2}}{n \sigma_{x}^{2}}} \qquad m_{r} = \sqrt{\frac{1 - r^{2}}{n - 2}}$$

- Определяются фактические значения t-критерия:
- Определяется критическое значение:

$$t_{\alpha;v_{OCT}}$$

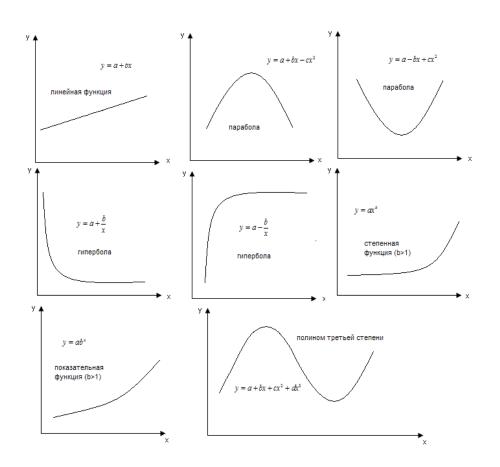
• Фактические значения сравниваются с критическими. Тестируемые параметры будут значимыми, если фактическое значение критерия превысит табличное





Способы определения формы уравнения регрессии:

- аналитический;
- графический;
- экспериментальный





Нелинейная регрессия

Два класса нелинейных регрессий

• регрессии, нелинейные относительно включенных в анализ объясняющих переменных, но линейные по оцениваемым параметрам. полиномы разных степеней

```
y = a + Bx + cx^2;

y = a + Bx + cx^2 + dx^3,

равносторонняя гипербола y = B + a/x.
```

• нелинейные регрессии по оцениваемым параметрам:

```
степенная y = a x^B показательная y = a B^X экспоненциальная y = e^{a+BX}
```





Подходы к линеаризации нелинейных уравнений

В параболе у = а + вх + сх², заменяя переменные x_1 =х, а x_2 =х², получаем двухфакторное уравнение линейной регрессии у = а + вх $_1$ + сх $_2$.

Для равносторонней гиперболы мы можем заменить 1/х на z и получим линейное уравнение регрессии, оценка параметров которого может быть дана МНК.

Для степенной функции у = а x^B применяется логарифмирование левой и правой части уравнения

 \log y = \log α + \log x $^{\beta}$, а затем замена переменных y_1 = \log y ; z = \log x и α_1 = \log α , получаем уравнение в линейном виде: y_1 = α_1 + β z



Спасибо за внимание!