

Лекция 10.

Свойство ковариации

Мет. Ковариацией случайных величин X и Y называется величина $\text{cov}(X, Y) = E((X - EX)(Y - EY)) = E(XY) - EX \cdot EY$

Ковариация является индикатором наличия направления связи между двумя случайными величинами

Пусть имеется $(X_1, Y_1), \dots, (X_n, Y_n)$ случайных величин X и Y

Def. Выборочной ковариацией называется величина $\widehat{\text{cov}}(X, Y) = \overline{xy} - \bar{x} \cdot \bar{y}$

По Закону Больших Чисел ясно, что $\widehat{\text{cov}}(X, Y) \rightarrow \text{cov}(X, Y)$, поэтому выборочная ковариация является оценкой

Th. Выборочная ковариация является точечной состоятельной, но смещенной оценкой ковариации. Несмещенной оценкой будет $\frac{n}{n-1} \widehat{\text{cov}}(X, Y)$

Ковариация и выборочная ковариация обладают свойствами

1. $\text{cov}(X, Y) = \text{cov}(Y, X)$
2. $\text{cov}(X, a) = 0$, где $a = \text{const}$
3. $\text{cov}(X, bY) = b \text{cov}(X, Y)$
4. $\text{cov}(X + Y, Z) = \text{cov}(X, Z) + \text{cov}(Y, Z)$
5. $\text{cov}(X, X) = D(X)$, $\widehat{\text{cov}}(X, X) = D^*(X)$
6. $D(X + Y) = DX + DY + 2 \text{cov}(X, Y)$

Nota. В дальнейшем под $\text{cov}(X, Y)$ будет пониматься выборочная ковариация

Анализ модели линейной парной регрессии

Пусть при n экспериментах получены значения случайных величин X и Z : $(X_1, Z_1), \dots, (X_n, Z_n)$

Пусть $X = \alpha + \beta Z + \varepsilon$ - теоретическая модель линейной регрессии, где ε - случайная величина, отражающая влияние невключенных факторов, нелинейность модели, ошибок измерений и просто случая.

Пусть построили с помощью метода наименьших квадратов выборочное уравнение линейной регрессии $\hat{X} = a + bZ$

Обозначим $\hat{\varepsilon}_i = X_i - \hat{X}_i + i$ - экспериментальная ошибка, разница между наблюдаемыми значениями и вычисляемыми по модели

Тогда $X_i = \hat{X}_i + \hat{\varepsilon}_i$ или $X_i = a + bZ_i + \hat{\varepsilon}_i$, где a и b - точечные оценки параметров α и β

Свойства $\hat{\varepsilon}_i$:

$$1. \overline{\hat{\varepsilon}} = 0$$

$$a = \bar{X} - b\bar{Z} \implies a + b\bar{Z} = \bar{X} \implies \overline{\hat{\varepsilon}} = \overline{X_i - (a + bZ_i)} = \bar{X} - \overline{a + bZ_i} = \bar{X} - \bar{X} = 0$$

$$2. \text{cov}(\hat{X}, \hat{\varepsilon}) = 0$$

$$\begin{aligned} b &= \overline{\bar{xz}} - \bar{x} \cdot \bar{z} \hat{\sigma}_z^2 = \overline{\text{cov}(X, Z)} D(Z) \implies \text{cov}(X, Z) - bD(Z) = 0 \\ \text{cov}(\hat{X}, \hat{\varepsilon}) &= \text{cov}(a + bZ, X - a - bZ) = \text{cov}(bZ, X - bZ) = \text{cov}(bZ, x) - \text{cov}(bZ, bZ) = \\ &= b \text{cov}(Z, X) - b^2 D(Z) = b(\text{cov}(Z, X) - bD(Z)) = 0 \end{aligned}$$

Анализ дисперсии результата

Def. $D(X) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ - дисперсия наблюдаемых значений

Def. $\hat{D}(X) = \frac{1}{n} \sum_{i=1}^n (\hat{X}_i - \bar{X})^2$ - дисперсия расчетных значений

Def. $D(\hat{\varepsilon}) = \frac{1}{n} \sum_{i=1}^n (\hat{\varepsilon}_i)^2$ - дисперсия остатков

Так как $X = \hat{X} + \hat{\varepsilon}$, $\text{cov}(\hat{X}, \hat{\varepsilon}) = 0$, то $D(X) = D(\hat{X}) + D(\hat{\varepsilon}) + 2 \text{cov}(\hat{X}, \hat{\varepsilon}) = D(\hat{X}) + D(\hat{\varepsilon})$

$$\text{Th. } D(X) = D(\hat{X}) + D(\hat{\varepsilon})$$

Очевидно, что качество модели будет тем лучше, чем меньше будет дисперсия остатков

Def. Коэффициентом детерминации R^2 называется величина $R^2 = \frac{D(\hat{X})}{D(X)}$ или $R^2 = 1 - \frac{D(\hat{\varepsilon})}{D(X)}$

Nota. Смысл R^2 - доля объясненной дисперсии, а $1 - R^2$ - доля необъясненной дисперсии

Свойства:

$$1. 0 \leq R^2 \leq 1$$

2. Если $R^2 = 1$, то $D(\hat{\varepsilon}) = 0 \implies \hat{\varepsilon}_i = \bar{\varepsilon} = 0$, то есть точки лежат строго на линии регрессии, модель идеальна

3. Если $R^2 = 0$, то $D(\hat{X}) = 0 \implies \hat{X} = \bar{x}$, то есть получаем примитивную, ничего не объясняющую модель

Чем больше R^2 , тем лучше качество модели

Проверка гипотезы о значимости уравнения регрессии

Проверяется $H_0 : R_{\text{теор}}^2 = 0$ (уравнение регрессии статистически не значимо) против $H_1 : R_{\text{теор}}^2 \neq 0$

Th. Если H_0 верна, то $F = \frac{R^2(n-2)}{1-R^2} \in F(1, n-2)$

Пусть t_α - квантиль $F(1, n-2)$ уровня α , тогда:

$$\begin{cases} H_0 : R_{\text{теор}}^2 = 0 & \text{если } F < t_\alpha \\ H_0 : R_{\text{теор}}^2 \neq 0 & \text{если } F \geq t_\alpha \end{cases}$$

Nota. Если $H_0 : R_{\text{теор}}^2 = 0$, то $H_0 : \beta = 0$

Связь между коэффициентом детерминации и коэффициентом линейной корреляции

1. $\sqrt{R^2} = r_{\hat{X}, X}$ - коэффициент корреляции между \hat{X} и X

$$r_{\hat{X}, X} = \frac{\text{cov}(\hat{X}, X)}{\sqrt{D(\hat{X})D(X)}} = \frac{\text{cov}(\hat{X}, \hat{X} + \varepsilon)}{\sqrt{D(\hat{X})D(X)}} = \frac{D(\hat{X}) + \cancel{\text{cov}(\hat{X}, \varepsilon)}^0}{\sqrt{D(\hat{X})D(X)}} = \sqrt{\frac{D(\hat{X})}{D(X)}} = R$$

2. $r_{\hat{X}, X} = |r_{X, Z}|$

$$\begin{aligned} \text{cov}(\hat{X}, X) &= \text{cov}(a + bZ, X) = b \text{cov}(Z, X) \\ D(\hat{X}) &= D(a + bZ) = b^2 D(Z) \\ r_{\hat{X}, X} &= \frac{\text{cov}(\hat{X}, X)}{\sqrt{D(\hat{X})D(X)}} = \frac{b \text{cov}(Z, X)}{\sqrt{b^2 D(Z)D(X)}} = \left| \frac{\text{cov}(X, Z)}{\sqrt{D(Z)D(X)}} \right| = |r_{X, Z}| \end{aligned}$$

Следствие 1: в случае линейной парной регрессии коэффициент детерминации равен квадрату коэффициенту корреляции

Следствие 2: в случае линейной парной регрессии совпадают результаты проверки гипотез

$$H_0 : R_{\text{теор}}^2 = 0 \iff H_0 : r = 0 \iff H_0 : \beta = 0$$

Теорема Гаусса-Маркова

Th. Пусть $X_i = \alpha + \beta Z_i + \varepsilon_i$ - теоретическая модель регрессии

$X = a + bZ$ - модель, полученная по методу наименьших квадратов

Если выполнено условия:

- а) Случайные члены ε_i независимые случайные величины, имеющие одинаковое нормальное распределение $\varepsilon_i \in N(0, \sigma^2)$
- б) Случайные величины ε_i и Z_i - независимы

Тогда a и b - состоятельные, несмещенные, эффективные оценки параметров α и β , то есть

1. Состоятельность: $a \xrightarrow[n \rightarrow \infty]{p} \alpha, b \xrightarrow[n \rightarrow \infty]{p} \beta$
2. Несмещенность: $Ea = \alpha, Eb = \beta$
3. Наименьшая дисперсия, равная:

$$Da = \frac{\overline{z^2} \sigma^2}{nD(Z)}, Db = \frac{\sigma^2}{nD(Z)}$$

Nota. Если не выполняется условие а), то есть ошибки зависимы или имеют разную дисперсию, то оценки становятся неэффективными. Если не выполнено условие б), то оценки становятся смещенными и несостоятельными

Стандартные ошибки коэффициентов регрессии

Из теоремы видим, что Da и Db зависят от дисперсии σ^2 случайного члена. По экспериментальным ошибкам получаем оценку данной дисперсии:

$$D(\hat{\varepsilon}) = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 \xrightarrow[n \rightarrow \infty]{p} \sigma^2$$

Однако эта оценка является смещенной:

$$E(D(\hat{\varepsilon})) = \frac{n-2}{n} \sigma^2$$

Поэтому несмещенной оценкой дисперсии σ^2 является величина $S^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2$

Def. Величина S называется стандартной ошибкой регрессии

Смысл: характеризует разброс наблюдаемых значений вокруг линии регрессии

Nota. Заменим в теореме Гаусса-Маркова σ^2 на S^2 , получаем оценки дисперсий Da и Db :

$$S_a^2 = \frac{\overline{z^2} S^2}{nD(z)}, S_b^2 = \frac{S^2}{nD(Z)}$$

Def. S_a и S_b называются стандартными ошибками коэффициентов регрессии

Прогнозирование регрессионных моделей

Пусть $X = \alpha + \beta Z + \varepsilon$ - теоретическая модель

$\hat{X} = a + bZ$ - модель МНК, построенная по выборке объема n

С помощью данной модели надо дать прогноз значения X_p при заданном значении Z_p и оценить качество прогноза

Теоретическое значение - $X_p = \alpha + \beta Z_p + \varepsilon$, а точечный прогноз $\hat{X}_p = a + bZ_p$

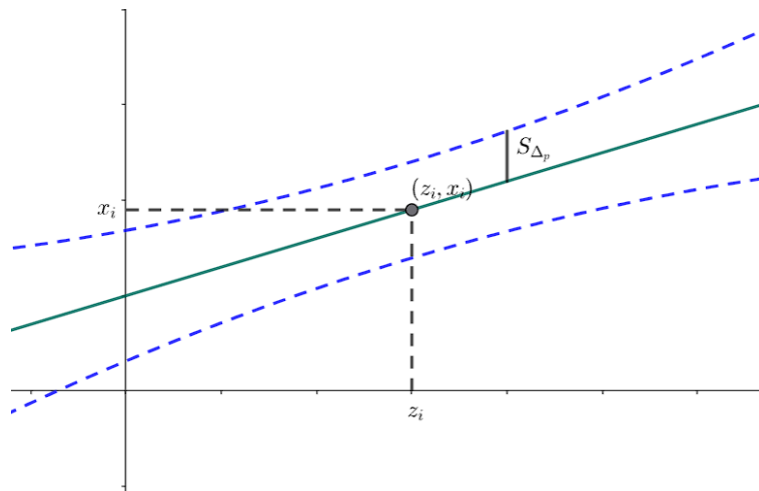
Разность между ними $\Delta_p = \hat{X}_p - X_p$ называется ошибкой предсказания

Свойства Δ_p :

1. $E\Delta_p = 0$
2. $D(\Delta_p) = \left(1 + \frac{1}{n} + \frac{(Z_p - \bar{z})^2}{nDZ}\right) \sigma^2$

Заменяя σ^2 на S^2 , получим стандартную ошибку прогноза: $S_{\Delta_p} = S \sqrt{1 + \frac{1}{n} + \frac{(Z_p - \bar{z})^2}{nDZ}}$

3. $D(\Delta_p) > \sigma^2$ - то есть точность прогноза ограничена случайным членом ε
4. При $n \rightarrow \infty$ $D(\Delta_p) \xrightarrow{p} \sigma^2$ - качество модели тем лучше, чем больше объем выборки
5. Чем больше Z_p отклоняется от \bar{z} , тем хуже качество прогноза. Наилучшее качество в точке $Z_p = \bar{z}$: $D(\Delta_p) = \left(1 + \frac{1}{n}\right) \sigma^2$



Доверительные интервалы прогноза и коэффициентов уравнения линейной регрессии

Пусть t_γ - квантиль $|T_{n-2}|$ уровня γ

Тогда доверительные интервалы надежности γ для параметров α и β :

$$\alpha : (a - t_\gamma S_a; a + t_\gamma S_a)$$

$$\beta : (b - t_\gamma S_b; b + t_\gamma S_b)$$

Доверительный интервал для прогноза X_p : $(\hat{X}_p - t_\gamma S_{\Delta_p}; \hat{X}_p + t_\gamma S_{\Delta_p})$