

# 1. Регрессионный анализ. Основы

## Основные задачи регрессионного анализа:

- а) подбираем класс функций для анализа;
- б) производим отбор наиболее информативных переменных;
- в) вычисляем оценки значений параметров модели;
- г) анализируем точность уравнения связи и его параметров;
- д) анализируем степень пригодности уравнения для целей прогноза.

# 1. Регрессионный анализ. Основы

Основные варианты статистических взаимосвязей между переменными  $X$  и  $Y$ :

- Не направленная связь, обе переменные равноценны, наличие и сила взаимосвязи – *корреляционный анализ*.
  - выделяем одну величину как *независимую* (*объясняющую, факторную, предиктор*), а другую как *зависимую* (*объясняемую, результирующую*), изменение первой может служить причиной для изменения второй переменной в виде какой либо **зависимости** – *регрессионный анализ*
- Количество, качество и неоднозначность связи.

# 1. Регрессионный анализ. Основы

Анализ влияния объясняющей переменной на зависимую переменную «в среднем»:

$$M(Y | x) = f(x)$$

*функция регрессии Y на X.* Здесь  $X$  – объясняющая переменная (*регрессор*),  $Y$  – зависимая.

Общий вид регрессии, суть:

$$M(Y | x_1, x_2, \dots, x_n) = f(x_1, x_2, \dots, x_n)$$

Парная, множественная регрессия. Ф. Далтон 19 в.

Возможность прогноза.

# 1. Регрессионный анализ. Основы

*Случайность* зависимости – не совпадение реальных значений зависимой переменной с её *условным математическим ожиданием*.

Дополнение случайным слагаемым  $\varepsilon$ , отражающим это несоответствие

$$Y = M(Y | x) + \varepsilon$$

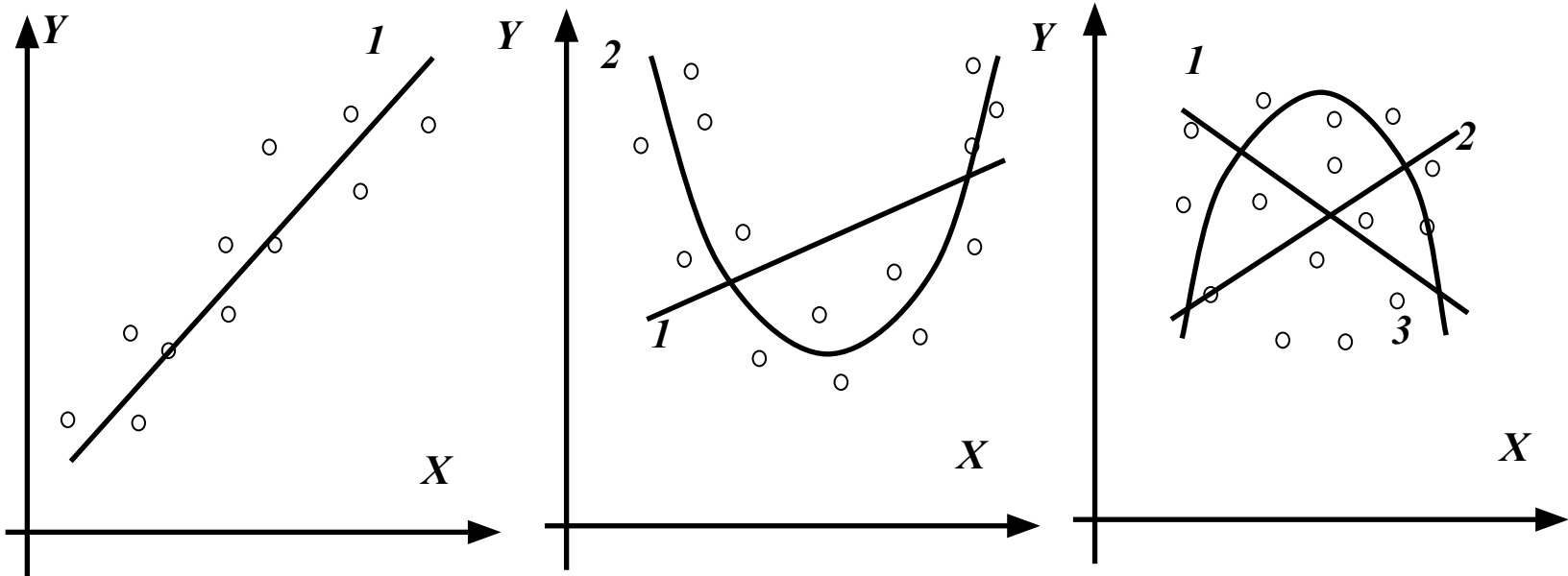
$$Y = M(Y | x_1, x_2, \dots, x_n) + \varepsilon$$

*Регрессионные модели (уравнения)*. Необходимость случайного члена.

# 1. Регрессионный анализ

1. Выбор вида связи переменных - спецификация уравнения *регрессии*.

Графический метод по *корреляционному полю* (*диаграмме рассеивания*). Аналитический.



# 1. Регрессионный анализ

2. Определение параметров (коэффициентов) модели - *параметризация*.

3. Проверка качества *регрессии* - *верификация*.

Линейная модель как общая тенденция процесса.

Параметризация линейной модели:

-Метод средних (Метод Маркова, на основе МЗР);

-Метод максимального правдоподобия.

Частный случай ММП при квадратичной целевой функции – метод наименьших квадратов МНК.

# 1. Регрессионный анализ

Задачи *линейного регрессионного анализа* для статистических данных  $(x_i, y_i)$  для переменных  $X$  и  $Y$ :

- а) получить наилучшие по какому-либо критерию оценки неизвестных теоретических параметров  $a_i$ ;
- б) проверить статистические гипотезы о параметрах модели;
- в) проверить адекватность модели данных наблюдений;
- г) выполнить оценку точности результатов вычисления
- д) провести анализ полученных данных.

# 1. Парный регрессионный анализ

Задача: подобрать для 2 рядов  $X$  и  $Y$  оптимальную модель из класса линейных вида

$$y_i = ax_i + b + e_i$$

чтобы наилучшим образом линейно выразить  $Y$  через  $X$ . С учетом коррекции имеем

$$\hat{y}_i = y_i + v_i = ax_i + b$$

Решение на основе:

- метода наименьших квадратов в натуральных величинах
- метода наименьших квадратов в центрированных величинах
- Байесовского метода (на основе характеристик многомерного закона распределения)



# 1. Парный регрессионный анализ

1. Для решения модели по МНК в натуральных величинах:

- на поправки  $v$  накладывают условие метода наименьших квадратов: найти такие значения коэффициентов  $a$  и  $b$  чтобы сумма квадратов поправок  $v - [v^2] = v^T v \rightarrow \min$  - была минимальной из всех возможных наборов пар коэффициентов.

Расчетные формулы:

- выражаем поправки  $v$ , записывая уравнения поправок для модели,

$$v_i = \hat{y}_i - y_i = ax_i + b - y_i$$

или в матричном виде

$$v = \hat{y} - y = Xk - y$$

Здесь матрица плана  $X$  и искомым коэффициентов имеет вид

$$k = \begin{pmatrix} a \\ b \end{pmatrix}$$

$$X = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \boxtimes & \boxtimes \\ x_n & 1 \end{pmatrix}$$

# 1. Парный регрессионный анализ

- составляем целевую функцию качества  $\Phi$  для МНК

$$\Phi = [v^2] = v^T v = [(ax + b - y)^2]$$

которая должна быть минимальна

- для минимизации берем частные производные от  $\Phi$  по  $k$  и полученное выражение приравниваем к нулю

$$\frac{\partial \Phi}{\partial k} = \frac{\partial \Phi}{\partial v} \cdot \frac{\partial v}{\partial k} = 2v^T \cdot X = 0$$

с учетом того, что  $v = Xk - y$

Транспонируя выражение имеем

$$X^T \cdot v = 0$$

- лемма Гаусса.

# 1. Парный регрессионный анализ

Подставим в лемму значение поправки  $v = Xk - y$

Получаем окончательно

$$X^T \cdot (X \cdot k - y) = X^T X \cdot k - X^T y = 0$$

матричную систему нормальных уравнений (совместная, 2 уравнения, 2 неизвестных) в свернутом виде

$$Nk = d,$$

с элементами

$$N = X^T X = \begin{bmatrix} N_{11} & N_{12} \\ N_{21} & N_{22} \end{bmatrix} = \begin{bmatrix} [x^2] & [x] \\ [x] & n \end{bmatrix}, \quad k = \begin{pmatrix} a \\ b \end{pmatrix}, \quad d = X^T y = \begin{pmatrix} d_1 \\ d_2 \end{pmatrix} = \begin{pmatrix} [xy] \\ [y] \end{pmatrix}$$

# 1. Парный регрессионный анализ

Систему в матричном виде целесообразно решать через обращение матрицы  $N$ , откуда имеем

$$k = \begin{pmatrix} a \\ b \end{pmatrix} = N^{-1} \cdot d = Q \cdot d$$

Здесь матрица  $Q = N^{-1}$  носит название обратная матрица.

Модельные значения в матричном виде будут

$$\hat{y} = XQX^T \cdot y = X \cdot k$$

Поправки в измерения получим как

$$v = \hat{y} - y = (XQX^T - E)y$$

# 1. Парный регрессионный анализ

2. Решение по МНК в центрированных величинах.

Для этого исходная модель

$$\hat{y}_i = y_i + v_i = ax_i + b$$

центрируется средним

$$\hat{y} - \bar{\hat{y}} = a(x - \bar{x}) \rightarrow \hat{y} = y + v = \bar{y} + a(x - \bar{x})$$

т.к.  $\hat{y} - \bar{\hat{y}} = \hat{y} - \bar{y}$

Имеем 1 неизвестное  $a$ . Тогда целевая функция  $\Phi$  будет

$$\Phi = [v^2] = [(\bar{y} + a(x - \bar{x}) - y)^2] = [(a \cdot v_x - v_y)^2]$$

Минимизация  $\Phi$  дает

$$\frac{\partial \Phi}{\partial a} = \frac{\partial \Phi}{\partial v} \cdot \frac{\partial v}{\partial a} = 2[(av_x - v_y) \cdot v_x] = 0$$

# 1. Парный регрессионный анализ

Откуда

- нормальное уравнение

$$\begin{bmatrix} v_x^2 \end{bmatrix} \cdot a = \begin{bmatrix} v_x \cdot v_y \end{bmatrix}$$

- решение относительно коэффициента  $a$

$$a = \frac{\begin{bmatrix} v_x \cdot v_y \end{bmatrix}}{\begin{bmatrix} v_x^2 \end{bmatrix}}$$

Для вычисления величины сдвига  $b$  раскроем скобки в центрированной модели модели и перегруппируем что получится

$$\hat{y} = a \cdot x + (\bar{y} - a \cdot \bar{x}) = a \cdot x + b$$

Откуда величина сдвига  $b$  будет

$$b = \bar{y} - a \cdot \bar{x}$$

# 1. Парный регрессионный анализ

Разделим числитель и знаменатель выражения для  $a$

$$a = \frac{[v_x \cdot v_y]}{[v_x^2]}$$

на число элементов в ряде  $n$ . Получим т.о. его оценку в терминах элементов ковариационной матрицы

$$a = \frac{\text{cov}(x,y)}{D(x)} = r_{xy} \cdot \frac{\sigma_y}{\sigma_x}$$

Подставив вид для  $a$  в центрированную модель, получим

$$\begin{aligned} \hat{y} &= \bar{y} + \frac{\text{cov}(x,y)}{D(x)} \cdot (x - \bar{x}) = \\ &= \bar{y} + r_{xy} \cdot \frac{\hat{\sigma}_y}{\hat{\sigma}_x} (x - \bar{x}) \end{aligned}$$

# 1. Парный регрессионный анализ

Полученное уравнение - оценка условного математического ожидания для нормально распределенной пары рядов. Метод получения уравнения регрессии в таком виде – метод средних, метод Маркова, метод условного математического ожидания, байесовский метод.

Этот же результат можно получить на основе выборочной ковариационной матрицы для 2 рядов  $x$  и  $y$ :

$$K = \begin{pmatrix} \frac{[v_x^2]}{n} & \frac{[v_x v_y]}{n} \\ \frac{[v_y v_x]}{n} & \frac{[v_y^2]}{n} \end{pmatrix} = \begin{pmatrix} D(x) & \text{cov}(x,y) \\ \text{cov}(y,x) & D(y) \end{pmatrix} = \begin{pmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{pmatrix}$$

Тогда имеем:



# 1. Парный регрессионный анализ

$$a = \frac{\text{cov}(y,x)}{D(x)} = \frac{K_{21}}{K_{11}}$$

$$b = (\bar{y} - a \cdot \bar{x}) = \bar{y} - \frac{K_{21}}{K_{11}} \cdot \bar{x}$$

$$\begin{aligned} \hat{y} &= \bar{y} + \frac{\text{cov}(x,y)}{D(x)} \cdot (x - \bar{x}) = \\ &= \bar{y} + K_{21} \cdot K_{11}^{-1} \cdot (x - \bar{x}) \end{aligned}$$

-т.е. все необходимые элементы модели в терминах ковариационной матрицы (2 мерный случай байесовского метода)

# 1. Парный регрессионный анализ

Если известна обратная ковариационная матрица  $K^{-1} = \Lambda$   
формулы для расчета коэффициентов примут вид

$$a = -\frac{\Lambda_{12}}{K_{22}}$$

и

$$b = (\bar{y} - a \cdot \bar{x}) = \bar{y} + \frac{\Lambda_{12}}{\Lambda_{22}} \cdot \bar{x}$$

Часто бывают более удобные чем другие.

В конце расчетов всегда записывается вид модели в числах.

# 1. Парный регрессионный анализ

Оценка точности уравнения регрессии:

- оценка модели в целом;
- оценка коэффициентов модели.

При необходимости – оценка значимости коэффициентов.

Чаще выполняется на основе метода наименьших квадратов и теоремы о переносе ошибок.

Для оценки точности модели:

- величины поправок  $v$  есть оценки истинных погрешностей  $\Delta$  модели;
- обычная формула стандартного отклонения для качества модели

$$\hat{\sigma}_0 = \sqrt{\frac{v^T v}{n-t}} = \sqrt{\frac{\Phi}{r}}$$

# 1. Парный регрессионный анализ

Получение погрешности модели на основе метода условного математического ожидания по условной дисперсии вида

$$\begin{aligned}\hat{\sigma}_{y|x}^2 &= \hat{\sigma}_0^2 = (K_{22} - K_{21} \cdot K_{11}^{-1} \cdot K_{12}) \cdot \left(\frac{n}{n-t}\right) = \left(K_{22} - \frac{K_{12}^2}{K_{11}}\right) \cdot \left(\frac{n}{n-t}\right) = \\ &= D(y) \cdot (1 - r_{xy}^2) \cdot \left(\frac{n}{n-t}\right)\end{aligned}$$

если используется ковариационная матрица и

$$\hat{\sigma}_0^2 = \frac{1}{\Lambda_{22}} \cdot \left(\frac{n}{n-t}\right)$$

Если использовать обратную ковариационную матрицу  $\Lambda$

# 1. Парный регрессионный анализ

Погрешности определения коэффициентов получим на основе теоремы о переносе ошибок используя формулу где вектор коэффициентов  $k$  линейно выражен через вектор измерений  $y$ :

$$k = QX^T \cdot y = F \cdot y$$

Тогда вид ковариационной матрицы  $K_k$  для вектора коэффициентов  $k$

$$K_k = F \cdot K_y \cdot F^T = (QX^T) \cdot K_y \cdot (XQ)$$

Вектор измерений  $y$  – один –  $K_y = \hat{\sigma}_0^2$  и

$$K_k = F \cdot K_y \cdot F^T = \hat{\sigma}_0^2 \cdot Q$$

размера  $(2 \times 2)$  - диагональные элементы – дисперсии вычисленных коэффициентов  $a$  и  $b$ .

# 1. Парный регрессионный анализ

Некоторые дополнительным возможностям регрессионного анализа:

- построение регрессии по методу наименьших квадратов для полиномиальной модели,
- представление на основе МНК периодических данных.

Очевидно, что это далеко не все дополнительные возможности парного регрессионного анализа.

Пусть появилось подозрение что исходные данные лучше будут представлены не линейной, а квадратичной моделью вида

$$\hat{y} = y + v = a \cdot x^2 + b \cdot x + c$$

или в матричном виде

$$y + v = X \cdot k$$

# 1. Парный регрессионный анализ

Целевая функция  $\Phi$  для модели при использовании МНК будет

$$\Phi = [v^2] = \left[ (a \cdot x^2 + b \cdot x + c - y)^2 \right]$$

производные от  $\Phi$  по коэффициентам  $a$ ,  $b$  и  $c$  приравненные к нулю

$$\frac{\partial \Phi}{\partial a} = \frac{\partial \left[ (a \cdot x^2 + b \cdot x + c - y)^2 \right]}{\partial a} = 2 \left[ (a \cdot x^2 + b \cdot x + c - y) \cdot x^2 \right] = 0$$

и т.д., откуда первое нормальное уравнение

$$\left[ x^4 \right] \cdot a + \left[ x^3 \right] \cdot b + \left[ x^2 \right] \cdot c = \left[ x^2 y \right]$$

и т.д., которых будет 3 по числу коэффициентов модели. Система совместная и разрешима относительно коэффициентов.

# 1. Парный регрессионный анализ

Рациональнее в матричном виде. Уравнения поправок

$$v = X \cdot k - y$$

где матрица плана  $X$  и свободный член  $y$  имеют вид

$$X = \begin{pmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ \boxtimes & \boxtimes & \boxtimes \\ x_n^2 & x_n & 1 \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \boxtimes \\ y_n \end{pmatrix}$$

целевая функция  $\Phi = v^T v$ , а её минимизация (производная по  $k$  и приравнивание к нулю дает лемму Гаусса

$$\frac{\partial \Phi}{\partial k} = \frac{\partial \Phi}{\partial v} \cdot \frac{\partial v}{\partial k} = 2v^T \cdot X = 0$$



# 1. Парный регрессионный анализ

Транспонируя выражение имеем более привычный вид

$$X^T \cdot v = 0$$

Подстановка в неё вида уравнений поправок  $v$  дает совместную систему нормальных уравнений

$$X^T \cdot (X \cdot k - y) = X^T X \cdot k - X^T y = 0$$

или

$$Nk = d,$$

с элементами

$$N = X^T X = \begin{pmatrix} N_{11} & N_{12} & N_{13} \\ N_{21} & N_{22} & N_{23} \\ N_{31} & N_{32} & N_{33} \end{pmatrix}, \quad k = \begin{pmatrix} a \\ b \\ c \end{pmatrix}, \quad d = X^T y = \begin{pmatrix} d_1 \\ d_2 \\ d_3 \end{pmatrix}$$

# 1. Парный регрессионный анализ

Подход может расширяться на любую степень полинома, изменяется только состав матрицы  $X$ .

Оценка точности стандартная для МНК:

- погрешность модели,
- погрешности коэффициентов,
- если необходимо, статистический анализ.

По погрешности модели можно сказать о необходимости усложнения модели.

# 1. Парный регрессионный анализ

Модель процесса может быть разная . Очень распространена периодическая функция представления данных.

При построении регрессии (аппроксимации) для периодической функции модель имеет следующий вид

$$\hat{y} = k_0 + k_1 \cdot \sin(x) + k_2 \cdot \cos(x) + k_3 \cdot \sin(2x) + k_4 \cdot \cos(2x) + \dots$$

Модели 1, 2 и т.д. порядков. Гармоники.

Для получения коэффициентов модели 1 порядка на основе МНК запишем уравнения поправок, целевую функцию  $\Phi$  и минимизируем её.

- уравнения поправок

$$v = k_0 + k_1 \cdot \sin(x) + k_2 \cdot \cos(x) - y$$

# 1. Парный регрессионный анализ

- целевая функция  $\Phi$

$$\Phi = [v^2] = \left[ (k_0 + k_1 \cdot \sin(x) + k_2 \cdot \cos(x) - y)^2 \right]$$

- Минимизация

$$\left\{ \begin{array}{l} \frac{\partial \Phi}{\partial k_0} = 2 \cdot [(k_0 + k_1 \cdot \sin(x) + k_2 \cdot \cos(x) - y) \cdot 1] = 0 \\ \frac{\partial \Phi}{\partial k_1} = 2 \cdot [(k_0 + k_1 \cdot \sin(x) + k_2 \cdot \cos(x) - y) \cdot \sin(x)] = 0 \\ \text{и т.д.} \end{array} \right.$$

и окончательно

$$\left\{ \begin{array}{l} n \cdot k_0 + [\sin(x)] \cdot k_1 + [\cos(x)] \cdot k_2 = [y] \\ [\sin(x)] \cdot k_0 + [\sin^2(x)] \cdot k_1 + [\cos(x) \cdot \sin(x)] \cdot k_2 = [y \cdot \sin(x)] \\ \dots \text{и т.д.} \end{array} \right.$$

# 1. Парный регрессионный анализ

Периодические функции имеют циклы и на основе свойств циклических перестановок имеем

$$[\sin(x)] = [\cos(x)] = [\cos(x) \cdot \sin(x)] = 0,$$

Откуда 
$$[\sin^2(x)] = [\cos^2(x)] = \frac{n}{2}$$

$$\left\{ \begin{array}{l} n \cdot k_0 + 0 \cdot k_1 + 0 \cdot k_2 = [y] \\ 0 \cdot k_0 + [\sin^2(x)] \cdot k_1 + 0 \cdot k_2 = [y \cdot \sin(x)] \\ \dots \text{и т.д.} \end{array} \right.$$

и явное решение

$$\left\{ \begin{array}{l} k_0 = \frac{[y]}{n} \\ k_1 = \frac{2 \cdot [y \cdot \sin(x)]}{n} \\ k_2 = \frac{2 \cdot [y \cdot \cos(x)]}{n} \end{array} \right.$$

# 1. Парный регрессионный анализ

Матрица системы нормальных уравнений имеет диагональный вид для модели любого порядка. Тогда по аналогии новые коэффициенты будут

$$\left\{ \begin{array}{l} k_3 = \frac{2 \cdot [y \cdot \sin(2x)]}{n} \\ k_4 = \frac{2 \cdot [y \cdot \cos(2x)]}{n} \\ \text{и т.д.} \end{array} \right.$$

а вычисленные ранее **не меняются**. Таким образом, модель 2 порядка получается из модели 1 порядка простым добавлением новых коэффициентов, 3 – добавление в модель 2 порядка вычисленных коэффициентов 3 порядка и т.д.

Оценка точности стандартная. Погрешность модели, погрешность коэффициентов, если надо – статистический анализ.

# 1. Парный регрессионный анализ

## Контрольные вопросы по теории 3 модуля

1. Основы общего регрессионного анализа. Основные этапы.
2. Парный линейный регрессионный анализ. Основные положения.
3. Решение задачи регрессии на основе МНК
4. Решение задачи регрессии на основе условного математического ожидания (метод средних).
5. Оценка точности и элементы статистического анализа в парной линейной регрессии.
6. Расширение парной модели регрессии на другие случаи.