

ТЕОРИЯ ВЕРОЯТНОСТЕЙ И МАТЕМАТИЧЕСКАЯ СТАТИСТИКА

Лекция 10

Изначальное содержание понятия регрессии:

Регрессия ([лат.](#) *regressio* — обратное движение, отход)

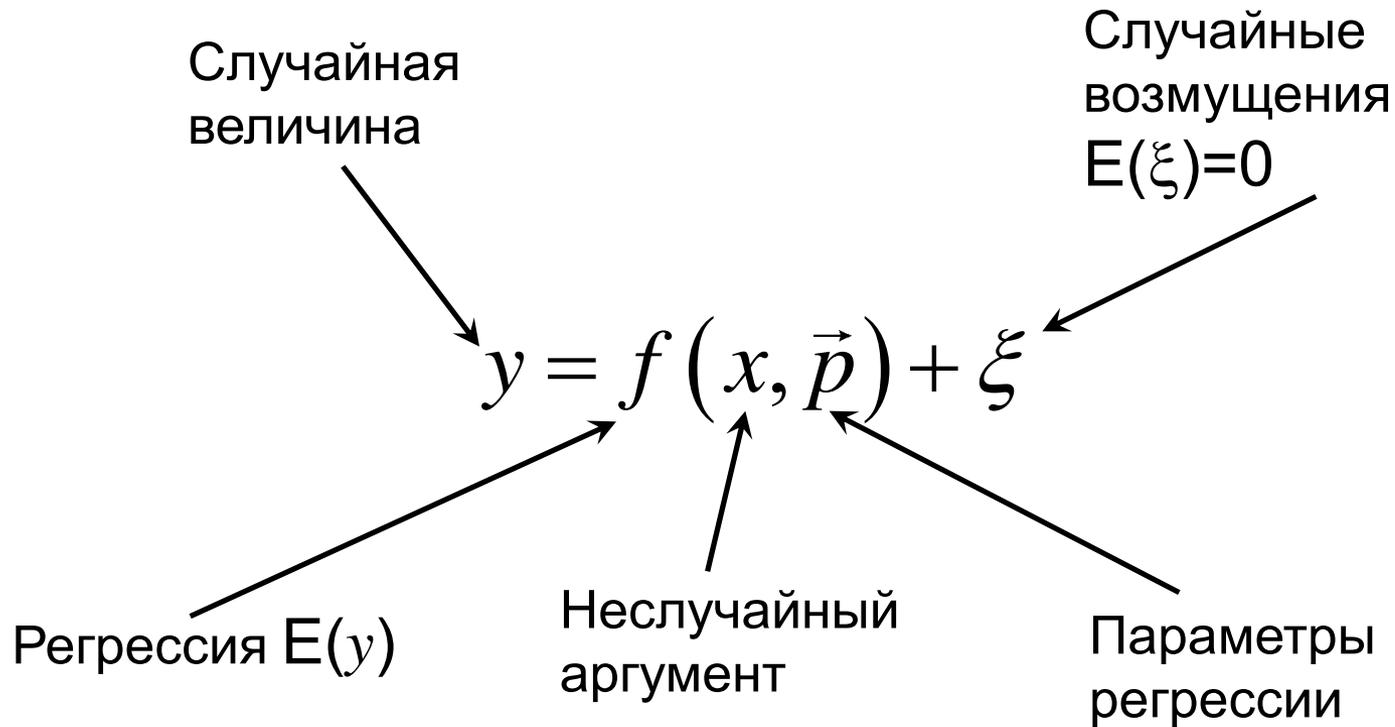
Общее понимание регрессии в математической статистике:

Пусть имеется пара случайных переменных (X, Y) .

Регрессия – это функция $g(x) = E(Y | X = x)$

то есть, условное математическое ожидание случайной переменной Y при условии, что другая случайная переменная X приняла значение x .

Далее регрессией будем называть зависимость математического ожидания случайной величины от неслучайных переменных.



Как правило, вид функции f задаётся исследователем.
Задача: найти значения вектора \vec{p} , наилучшим образом соответствующие данным, полученным в эксперименте, когда значения x задаются, а значения y измеряются.

Если f - полином, то

$$f(x, \vec{p}) = a_0 + a_1x + a_2x^2 + a_3x^3 + \dots + a_nx^n$$

где n - конечное и небольшое.

Вектор параметров:

$$\vec{p} = (a_0, a_1, a_2, a_3, \dots, a_n)$$

Это пример линейной регрессии.

Общий случай линейной регрессии:

$$f(x, \vec{p}) = \sum_{l=1}^L p_l f_l(x)$$

f_l - базисные функции, не зависящие от p_l

В случае полинома

$$f_1(x) = x^0 = 1, \quad f_2(x) = x^1 = x, \quad f_3(x) = x^2 \quad \dots \quad f_L(x) = x^{L-1}$$

Вектор параметров \vec{p} должен быть таким, чтобы кривая регрессии наилучшим образом прошла между экспериментальными точками y

Принцип максимального правдоподобия

Предполагается, что $w(\xi) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\xi^2}{2\sigma^2}}$

Отсюда $w(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{[y-f(x,\vec{p})]^2}{2\sigma^2}}$

(x_n, y_n) – пары экспериментальных данных

Вероятность попадания y_n в интервал $y_n - \frac{1}{2}\Delta y < y < y_n + \frac{1}{2}\Delta y$:

$$w(y_n) = \frac{\Delta y}{\sqrt{2\pi\sigma^2}} e^{-\frac{[y_n - f(x_n, \vec{p})]^2}{2\sigma^2}}$$

Для N экспериментальных точек:

$$w(y_1, y_2, \dots, y_N) = \prod_{n=1}^N w(y_n) = \left(\frac{\Delta y}{\sqrt{2\pi\sigma^2}} \right)^N \prod_{n=1}^N e^{-\frac{[y_n - f(x_n, \vec{p})]^2}{2\sigma^2}}$$

Функция правдоподобия:

$$LH = \prod_{n=1}^N e^{-\frac{[y_n - f(x_n, \vec{p})]^2}{2\sigma^2}} = \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N [y_n - f(x_n, \vec{p})]^2 \right\}$$

Функция правдоподобия максимальна при

$$S = \sum_{n=1}^N [y_n - f(x_n, \vec{p})]^2 = \min$$

На этом основан метод наименьших квадратов:
найти значения параметров \vec{p} ,
при которых S минимальна

Метод наименьших квадратов:

$$\frac{\partial}{\partial p_l} S = -2 \sum_{n=1}^N \left\{ \left[y_n - f(x_n, \vec{p}) \right] \frac{\partial}{\partial p_l} f(x_n, \vec{p}) \right\} = 0, \quad l = 1, 2, \dots, L$$

Общий случай
линейной регрессии:

$$f(x, \vec{p}) = \sum_{l=1}^L p_l f_l(x)$$

Тогда

$$\frac{\partial}{\partial p_l} f(x, \vec{p}) = f_l(x)$$

Подставляем:

$$\sum_{n=1}^N \left\{ \left[y_n - \sum_{m=1}^L p_m f_m(x_n) \right] f_l(x_n) \right\} = 0, \quad l = 1, 2, \dots, L$$

$$\sum_{n=1}^N y_n f_l(x_n) - \sum_{n=1}^N \left\{ f_l(x_n) \sum_{m=1}^L p_m f_m(x_n) \right\} = 0, \quad l = 1, 2, \dots, L$$

$$\sum_{m=1}^L p_m \left\{ \sum_{n=1}^N f_l(x_n) f_m(x_n) \right\} = \sum_{n=1}^N y_n f_l(x_n), \quad l = 1, 2, \dots, L$$

Обозначим:

$$\sum_{n=1}^N f_l(x_n) f_m(x_n) = a_{lm}, \quad \sum_{n=1}^N y_n f_l(x_n) = b_l$$

Система линейных уравнений
для параметров линейной регрессии:

$$\sum_{m=1}^L a_{lm} p_m = b_l, \quad l = 1, 2, \dots, L$$

Векторно-матричная форма:

$$A\vec{p} = \vec{b}$$

Решение:

$$\vec{p} = A^{-1}\vec{b}$$

Расчет матрицы ошибок значений параметров:

$$RSS = S_{\min} \quad s_y^2 = \frac{RSS}{N - L}$$

Матрица ошибок:

$$D_{\vec{p}} = s_y^2 A^{-1}$$

$$D_{\vec{p}} = \begin{pmatrix} \sigma_{p_1}^2 & \text{cov}(p_1, p_2) & \dots & \text{cov}(p_1, p_L) \\ \text{cov}(p_2, p_1) & \sigma_{p_2}^2 & \dots & \text{cov}(p_2, p_L) \\ \dots & \dots & \dots & \dots \\ \text{cov}(p_L, p_1) & \text{cov}(p_L, p_2) & \dots & \sigma_{p_L}^2 \end{pmatrix}$$

$$r(p_i, p_j) = \frac{\text{cov}(p_i, p_j)}{\sigma_{p_i} \sigma_{p_j}}$$

Алгоритм вычисления параметров линейной регрессии

1) Выбор базисных функций

Затем вычисление следующих величин

- 2) Матрица A^{-1} и вектор правой частей \vec{b}
- 3) Матрица \vec{p}
- 4) Параметры регрессии
- 5) Остаточная сумма квадратов отклонений RSS и величина \sum_y^2
- 6) Матрица ошибок $D_{\vec{p}}$