

Множественная линейная регрессия. Постановка задачи. Вероятностный подход. Решение задачи. Прогнозирование.

Напомним, что *регрессионный анализ* может одновременно рассматриваться и как раздел *математической статистики*, и как раздел *эконометрики*; *регрессионному анализу* посвящена глава в большинстве монографий и учебников по *анализу данных*.

Постановка задачи.

Предположим, что имеется n объектов, каждый из которых описывается m признаками (*факторами, предикторами*), влияющими на значение *переменной отклика*. Линейная регрессионная модель необходима нам для прогноза значений переменной отклика по известным значениям факторов.

Например, весьма актуально получение инструмента для оценки возможной *стоимости квартиры* по её основным характеристикам, таким как

- общая площадь;
- жилая площадь;
- расстояние до центра;
- расстояние до метро;
- этаж;
- этажность дома
- ...

В терминах регрессионной модели,

- под *объектами* будем понимать квартиры;
- под *факторами* – перечисленные выше свойства;
- под *переменной отклика* – стоимость квартиры.

Будем нумеровать объекты индексом i ($i = 1, \dots, n$), а факторы – индексом j ($j = 1, \dots, m$).

Обозначим через x_{ij} – значение j -го признака i -го объекта ($i = 1, \dots, n$; $j = 1, \dots, m$). Таким образом, каждый i -ый объект представим как m -мерный вектор $X_i \equiv (x_{i1}, \dots, x_{im})$, $i = 1, \dots, n$.

Сырые данные для этой задачи могут иметь, например, следующий вид:

TotalSquare (m2)	LivingSquare (m2)	DistCenter (km)	DistMetro (km)	Price
80	53	17	2,1	14 612 000,00 ₺
76	51	1	0,7	16 931 128,00 ₺
96	72	16	1,3	18 905 472,00 ₺
56	37	16	2,2	14 829 304,00 ₺
75	56	6	2,8	19 214 025,00 ₺
75	56	11	1,1	19 582 950,00 ₺
97	65	12	1,4	19 123 259,00 ₺
30	24	14	2,3	6 035 280,00 ₺
84	63	7	1,9	20 058 696,00 ₺
50	33	11	2,4	13 807 800,00 ₺
55	44	6	1,7	13 087 745,00 ₺
94	71	10	1,7	17 337 266,00 ₺
91	68	5	1,4	17 189 900,00 ₺
32	26	7	1,8	6 405 792,00 ₺
86	65	4	0,3	19 267 698,00 ₺
55	41	2	1,6	13 827 495,00 ₺
65	49	18	1,7	11 242 920,00 ₺
45	34	9	2,2	12 004 470,00 ₺
47	38	4	1,3	10 586 844,00 ₺

Здесь каждая строка соответствует одной квартире.

Обозначим значение *зависимого признака* i -го объекта (т.е. стоимости квартиры) через y_i , $i = 1, \dots, n$.

Будем искать зависимость в виде линейной функции m переменных:

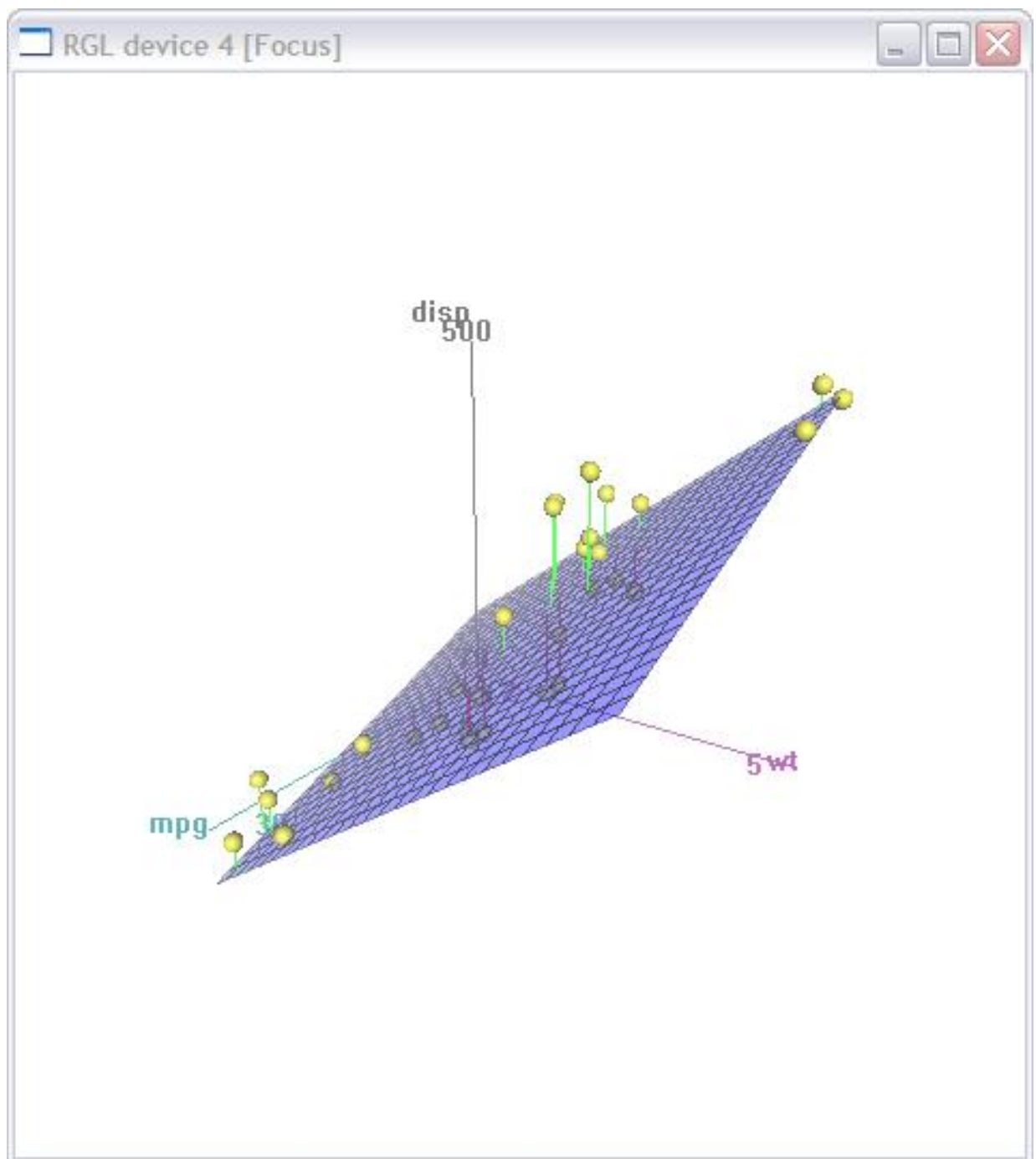
$$f(x_1, \dots, x_m) = \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m, \quad (1)$$

Здесь x_j , ($j = 1, \dots, m$) – независимые переменные (факторы, т.е. характеристики квартиры), зависимая переменная (прогнозируемая цена квартиры), $\theta_0, \theta_1, \dots, \theta_m$ – *искомые параметры* линейной зависимости. При $m=1$ получаем рассмотренную ранее модель *однофакторной линейной регрессии*.

Геометрическая интерпретация

Геометрически можно интерпретировать задачу следующим образом. Точки с координатами $(x_{i1}, \dots, x_{im}, y_i)$, $i = 1, \dots, n$, образуют «облако» в пространстве R^{m+1} . Нужно построить такую *гиперплоскость* вида (1),

которая бы наилучшим образом проходила через это облако точек. Здесь термин «наилучший» означает (как и в случае однофакторной линейной регрессии) «минимизирующий сумму квадратов отклонений фактических значений зависимой переменной от ожидаемых в соответствии с линейной моделью». Для случая 2-х переменных линейная модель представима как плоскость в 3-хмерном пространстве:



Источник: <http://www.statmethods.net/graphs/images/scatter3d.png>

Решение задачи. Вероятностный подход.

Понятно, что вследствие наличия случайной компоненты точки «облака», вообще говоря, не будут находиться в одной гиперплоскости.

Обозначим случайную величину («ошибку») i -го измерения через ε_i . Будем считать, что ошибки всех n измерений распределены одинаково, а именно – закон распределения вероятности ошибки будем предполагать *нормальным (гауссовым)* и будем считать, что *систематическая* ошибка отсутствует (то есть «разброс» цен при одинаковых параметрах квартир вызван исключительно совокупным влиянием случайных факторов). В этом случае математическое ожидание ошибки равно нулю. Обозначим среднее квадратическое отклонение ошибки через σ . Математическая запись названных условий имеет вид: $M\varepsilon_i = 0, D\varepsilon_i = \sigma^2, i = 1, \dots, n$).

Запишем систему относительно искомых величин $\theta_0, \theta_1, \dots, \theta_m$:

$$\begin{aligned} y_1 &= \theta_0 + \theta_1 x_{11} + \dots + \theta_m x_{1m} + \varepsilon_1 \\ y_2 &= \theta_0 + \theta_1 x_{21} + \dots + \theta_m x_{2m} + \varepsilon_2 \\ &\dots \dots \dots \\ y_i &= \theta_0 + \theta_1 x_{i1} + \dots + \theta_m x_{im} + \varepsilon_i, \\ &\dots \dots \dots \\ y_n &= \theta_0 + \theta_1 x_{n1} + \dots + \theta_m x_{nm} + \varepsilon_n \end{aligned} \quad (2)$$

Введём обозначения:

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1m} \\ 1 & x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2m} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{im} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{nm} \end{pmatrix} \quad Y = \begin{pmatrix} y_1 \\ \dots \\ y_n \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \dots \\ \varepsilon_n \end{pmatrix}, \quad \Theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \dots \\ \theta_m \end{pmatrix}$$

и перепишем систему (2) в матричной форме:

$$Y = X\Theta + \varepsilon. \quad (3)$$

Выразим вектор ошибок \mathcal{E} из системы (3):

$$\mathcal{E} = Y - X\Theta. \quad (4)$$

Запишем сумму квадратов ошибок (минимизируемую функцию):

$$\begin{aligned} \sum_{i=1}^n \mathcal{E}_i^2 &\equiv \mathcal{E}^T \mathcal{E} = (Y - X\Theta)^T (Y - X\Theta) = (Y^T - \Theta^T X^T)(Y - X\Theta) = \\ &= Y^T Y - Y^T X\Theta - \Theta^T X^T Y + \Theta^T X^T X\Theta = Y^T Y - 2\Theta^T X^T Y + \Theta^T X^T X\Theta. \end{aligned}$$

Обозначим

$$F(\Theta) \equiv Y^T Y - 2\Theta^T X^T Y + \Theta^T X^T X\Theta. \quad (5)$$

Минимизируем функцию (5) по Θ . Найдём градиент функции $F(\Theta)$:

$\text{grad } F(\Theta) = -2X^T Y + 2X^T X\Theta$ и потребуем выполнения равенства $\text{grad } F(\Theta) = 0$. Получим систему:

$$X^T X\Theta = X^T Y. \quad (6)$$

Будем считать, что $\det(X^T X) \neq 0$, а значит, решение системы (6) запишется в виде:

$$\hat{\Theta} = (X^T X)^{-1} X^T Y. \quad (7)$$

Найденное решение (вектор параметров) позволяет осуществлять прогноз значений зависимого признака по известным значениям независимых признаков:

Применительно к рассмотренному примеру, получим прогнозируемое значение

$$\hat{y}_i = X_i \hat{\Theta}, \quad i = n+1, \dots \quad (8)$$