

# Спецкурс "Регрессия"

Шкляев А.В.

21 января 2019 г.

## 1 Нормальная модель

### 1.1 Некоторые замечания

На прошлой лекции мы получили выражения для оценок в общей линейной модели  $\vec{Y} = X\vec{a} + \vec{\varepsilon}$ , где

$$\hat{a} = (X^t X)^{-1} X^t \vec{Y}, \quad \hat{\sigma}^2 = \frac{RSS}{n} = \frac{\|\vec{Y} - X\hat{a}\|^2}{n}.$$

Отметим также, что

$$RSS = \|\vec{Y}(E - X(X^t X)^{-1} X^t)\|^2 = \vec{Y}^t (E - X(X^t X)^{-1} X^t)^t (E - X(X^t X)^{-1} X^t) \vec{Y} = \\ \vec{Y}^t (E - 2X(X^t X)^{-1} X^t + X(X^t X)^{-1} X^t) \vec{Y} = \vec{Y}^t (E - X(X^t X)^{-1} X^t) \vec{Y}.$$

Также как и прежде можно оценить  $\mathbf{E}(Y^*|X^*) = \langle \vec{X}^*, \vec{a} \rangle$ . При этом

$$(\vec{X}^*)^t (X^t X)^{-1} X^t \vec{Y} = (\vec{X}^*)^t \vec{a} + (\vec{X}^*)^t (X^t X)^{-1} X^t \vec{\varepsilon},$$

откуда

$$\hat{Y}^* \sim \mathcal{N}(\langle \vec{X}^*, \vec{a} \rangle, \sigma^2 (E + \Sigma)), \quad \Sigma^2 = (\vec{X}^*)^t (X^t X)^{-1} X^t X (X^t X)^{-1} \vec{X}^* = (\vec{X}^*)^t (X^t X)^{-1} \vec{X}^*.$$

Следовательно,  $\hat{Y}^*$  несмещенная оценка  $\mathbf{E}(Y^*|X^*)$ .

Сделаем также следующее замечание. Если матрица  $X$  неполного ранга (в нашей задаче это означает, что столбцы матрицы линейно зависимы), матрица  $X^t X$  необратима и указанные выше формулы для подсчета коэффициентов неприменимы. В этом случае, как нетрудно понять, решение  $\hat{a}$  не единственно, а определяется соотношением

$$X^t \vec{y} = X^t X \hat{a}. \quad (1)$$

При этом мы по-прежнему находим проекцию  $\vec{Y}$  на пространство  $L = \{X\vec{a}, \vec{a} \in \mathbb{R}^n\}$ , просто пространство  $L$  плохо (избыточно) параметризовано. Мы не можем найти оценки параметров  $\vec{a}$ , обладающие сколько-то хорошими свойствами, но при любом выборе  $\vec{a}$  мы получим одно и то же RSS — квадрат расстояния между  $\vec{y}$  и  $L$ . Следовательно, оценка  $\sigma^2$  будет обладать теми же свойствами, что и прежде.

### 1.2 Выбор матрицы плана

Предположим, что мы можем выбрать матрицу  $X$ . Как выбрать ее таким образом, чтобы снизить дисперсии оценок  $\hat{a}_i$ ? Мы знаем явные выражения  $\mathbf{D}\hat{a}_i = (X^t X)^{-1}_{i,i} \sigma^2$ .

Минимизируем диагональные элементы, наложив ограничения на столбцы матрицы  $X$  вида  $\|x_{\cdot,i}\| \leq c_i$ .

При этом

$$A = X^t X = \begin{pmatrix} \|x_{\cdot,1}\|^2 & \vec{b}^t \\ \vec{b} & F \end{pmatrix}, \quad F = \tilde{X}^t \tilde{X}, \quad \tilde{X}_{i,j} = \langle x_{\cdot,i+1}, x_{\cdot,j+1} \rangle, \quad \vec{b}_i = \langle x_{\cdot,1}, x_{\cdot,i+1} \rangle.$$

Следовательно,

$$\mathbf{D}\hat{a}_1 = A_{i,i}^{-1} = \det F / \det A.$$

Свяжем  $\det(F)$  и  $\det(A)$ , используя матрицу

$$C = \begin{pmatrix} 1 & 0 \\ -F^{-1}\vec{b} & E \end{pmatrix}, \quad AC = \begin{pmatrix} \|x_{\cdot,1}\|^2 & \vec{b}^t \\ \vec{b} & F \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -F^{-1}\vec{b} & E \end{pmatrix} = \begin{pmatrix} \|x_{\cdot,1}\|^2 - \vec{b}^t F^{-1} \vec{b} & \vec{b}^t \\ 0 & F \end{pmatrix}.$$

Поскольку  $\det(C) = 1$ , то

$$\det(A) = \det(AC) = \left( \|x_{\cdot,1}\|^2 - \vec{b}^t F^{-1} \vec{b} \right) \det(F).$$

Отсюда

$$\mathbf{D}\hat{a}_1 = \frac{1}{\|x_{\cdot,1}\|^2 - \vec{b}^t F^{-1} \vec{b}}$$

будет минимальной при минимальной  $\vec{b}^t F^{-1} \vec{b}$ . При этом матрица  $F$  является положительно определенной, поскольку

$$\vec{u}^t F \vec{u} = \vec{u}^t \tilde{X}^t \tilde{X} \vec{u} = \langle \tilde{X} \vec{u}, \tilde{X} \vec{u} \rangle > 0$$

при  $\vec{u} \neq 0$  (так как  $\tilde{X}$  полного ранга). Следовательно,  $F^{-1}$  также положительно определена и  $\vec{b}^t F^{-1} \vec{b} \geq 0$ , где равенство возможно только при  $\vec{b} = 0$ . Значит минимум  $\mathbf{D}\hat{a}_{i,i}$  достигается при максимальной доступной норме  $\|x_{\cdot,1}\| = c_1$  и при  $\vec{b} = 0$ , что в силу определения  $\vec{b}$  означает ортогональность  $x_{\cdot,1}$  и  $x_{\cdot,i}$  при  $i > 1$ .

Аналогичные рассуждения приводят нас к выводу о том, что все  $x_{\cdot,i}$  должны быть ортогональными и иметь максимальную возможную норму.

### 1.3 Общая модель без предположения нормальности

Напомню, что идеологически мы рассматриваем задачу

$$\mathbf{E}(Y - f(\vec{X}))^2 \rightarrow \min \quad (2)$$

в классе функций  $f(\vec{u}) = \langle \vec{u}, \vec{a} \rangle$ . В случае, если мы не знаем распределения  $(X, Y)$ , но знаем, что  $Y_i = \varepsilon_i + \langle X_{i,\cdot}, \vec{a} \rangle$ , где  $\varepsilon_i$  н.о.р., то естественно заменить задачу о минимизации (2) на минимизацию

$$\sum_{i=1}^n (Y_i - \langle X_{i,\cdot}, \vec{a} \rangle)^2.$$

Вторая задача приводит нас к той же оценке МНК, что и в нормальной модели. При этом решение задачи минимизации  $\|\vec{y} - X\vec{a}\|$  не связано с распределениями (да и вообще с теорией вероятностей), поэтому решением ее как и прежде будет

$$\hat{a} = (X^t X)^{-1} X^t \vec{Y}.$$

Более того, это как и прежде будет несмещенная оценка  $\vec{a}$  с ковариационной матрицей  $\sigma^2(X^t X)^{-1}$ . Аналогичным образом

$$\hat{\sigma}^2 = RSS/(n-2), \quad RSS = \|\vec{Y} - X\hat{a}\|^2,$$

будет несмещенной оценкой  $\sigma^2$ .

Такая оценка называется оценкой методом наименьших квадратов или МНК-оценкой. Как мы видели, в нормальной модели эта оценка была оценкой максимального правдоподобия. В общей модели это не так, но можно сформулировать следующую теорему:

**Теорема 1** (Теорема Гаусса-Маркова). Пусть  $\vec{c}$  — детерминированный вектор. Тогда величина  $\langle \vec{c}, \hat{a} \rangle$  есть несмещенная оценка  $\langle \vec{a}, \vec{a} \rangle$  с наименьшей дисперсией среди всех несмещенных оценок вида  $d_1(X)Y_1 + \dots + d_n(X)Y_n$ .

*Доказательство.* Как и прежде будем рассматривать фиксированные  $x_{i,j}$ . Наша оценка имеет вид

$$\langle \vec{c}, \hat{a} \rangle = \langle \vec{c}, (X^t X)^{-1} X^t \vec{Y} \rangle = \vec{Y}^t X (X^t X)^{-1} \vec{c}.$$

Для произвольной несмещенной оценки вектор  $\vec{d}(X) = (d_1(X), \dots, d_n(X))$  удовлетворяет условиям

$$\mathbf{E} \left( \sum_{i=1}^n d_i(X) Y_i \right) = \vec{d}(X)^t X \vec{a} = \vec{c}^t \vec{a}$$

при всех  $\vec{a}$ , то есть  $X^t \vec{d}(X) = \vec{c}$ . Среди таких векторов мы должны выбрать вектор минимальной длины  $\|\vec{d}\|$ , поскольку

$$\mathbf{D} \left( \sum_{i=1}^n d_i(X) Y_i \right) = \sigma^2 \sum_{i=1}^n d_i(X)^2.$$

Минимальную длину среди таких  $\vec{d}$  имеет перпендикуляр из нуля на плоскость  $X^t \vec{d} = \vec{c}$ . Убедимся, что вектор  $\vec{d}_0 = X (X^t X)^{-1} \vec{c}$  ортогонален всем векторам внутри этой плоскости. Векторы внутри плоскости удовлетворяют соотношению  $X^t \vec{d} = 0$ , при этом

$$\langle \vec{d}_0, \vec{d} \rangle = \vec{c}^t (X^t X)^{-1} X^t \vec{d} = 0.$$

Итак,  $\vec{d}_0$  ортогонален плоскости  $X^t \vec{d} = \vec{c}$  и лежит в ней. Значит он имеет минимальную длину из векторов этой плоскости, откуда и вытекает теорема Гаусса-Маркова.  $\square$