

$f_1(x), \dots, f_n(x)$ — числовые признаки;

Модель многомерной линейной регрессии:

$$f(x, \alpha) = \sum_{j=1}^n \alpha_j f_j(x), \quad \alpha \in \mathbb{R}^n.$$

Матричные обозначения:

$$F_{\ell \times n} = \begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}, \quad y_{\ell \times 1} = \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}, \quad \alpha_{n \times 1} = \begin{pmatrix} \alpha_1 \\ \dots \\ \alpha_n \end{pmatrix}.$$

Функционал квадрата ошибки:

$$Q(\alpha, X^\ell) = \sum_{i=1}^{\ell} (f(x_i, \alpha) - y_i)^2 = \|F\alpha - y\|^2 \rightarrow \min_{\alpha}.$$

Необходимое условие минимума в матричном виде:

$$\frac{\partial Q}{\partial \alpha}(\alpha) = 2F^T(F\alpha - y) = 0,$$

откуда следует *нормальная система* задачи МНК:

$$F^T F \alpha = F^T y,$$

где $F^T F$ — $n \times n$ -матрица.

Решение системы: $\alpha^* = (F^T F)^{-1} F^T y = F^+ y$.

Значение функционала: $Q(\alpha^*) = \|P_F y - y\|^2$,

где $P_F = FF^+ = F(F^T F)^{-1} F^T$ — проекционная матрица.

Произвольная $\ell \times n$ -матрица представима в виде сингулярного разложения (singular value decomposition, SVD):

$$F = VDU^T.$$

Основные свойства сингулярного разложения:

- ❶ $\ell \times n$ -матрица $V = (v_1, \dots, v_n)$ ортогональна, $V^T V = I_n$, столбцы v_j — собственные векторы матрицы FF^T ;
- ❷ $n \times n$ -матрица $U = (u_1, \dots, u_n)$ ортогональна, $U^T U = I_n$, столбцы u_j — собственные векторы матрицы $F^T F$;
- ❸ $n \times n$ -матрица D диагональна, $D = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$, $\lambda_j \geq 0$ — собственные значения матриц $F^T F$ и FF^T , $\sqrt{\lambda_j}$ — сингулярные числа матрицы F .

Псевдообратная F^+ , вектор МНК-решения α^* ,
МНК-аппроксимация целевого вектора $F\alpha^*$:

$$F^+ = (UDV^T VDU^T)^{-1}UDV^T = UD^{-1}V^T = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j v_j^T;$$

$$\alpha^* = F^+ y = UD^{-1}V^T y = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y);$$

$$F\alpha^* = P_F y = (VDU^T)UD^{-1}V^T y = VV^T y = \sum_{j=1}^n v_j (v_j^T y);$$

$$\|\alpha^*\|^2 = \|D^{-1}V^T y\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^T y)^2.$$

Проблема: мультиколлинеарность при $\lambda_j \rightarrow 0$.

Если имеются сингулярные числа, близкие к нулю, то:

- матрица $\Sigma = F^T F$ плохо обусловлена;
- решение становится неустойчивым и неинтерпретируемым, слишком большие коэффициенты $\|\alpha_j^*\|$ разных знаков;
- возникает переобучение:
на обучении $Q(\alpha^*, X^\ell) = \|F\alpha^* - y\|^2$ мало;
на контроле $Q(\alpha^*, X^k) = \|F'\alpha^* - y'\|^2$ велико;

Стратегии устранения мультиколлинеарности и переобучения:

- отбор признаков: $f_1, \dots, f_n \rightarrow f_{j_1}, \dots, f_{j_m}, \quad m \ll n$.
- регуляризация: $\|\alpha\| \rightarrow \min$;
- преобразование признаков: $f_1, \dots, f_n \rightarrow g_1, \dots, g_m, \quad m \ll n$;

- Задача многомерной линейной регрессии может быть решена через сингулярное разложение
- Мультиколлинеарность приводит к плохой обусловленности, неустойчивости и переобучению
- Методы устранения мультиколлинеарности (гребневая регрессия, метод главных компонент) также связаны с сингулярным разложением (об этом в следующих лекциях)