

# Многомерная линейная регрессия

$f_1(x), \dots, f_n(x)$  — 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:

- 1  $\ell \times n$ -матрица  $V = (v_1, \dots, v_n)$  ортогональна,  $V^T V = I_n$ , столбцы  $v_j$  — собственные векторы матрицы  $FF^T$ ;
- 2  $n \times n$ -матрица  $U = (u_1, \dots, u_n)$  ортогональна,  $U^T U = I_n$ , столбцы  $u_j$  — собственные векторы матрицы  $F^T F$ ;
- 3  $n \times n$ -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  $\alpha^*$ ,  
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)

## Регуляризация (гребневая регрессия)

Штраф за увеличение нормы вектора весов  $\|\alpha\|$ :

$$Q_\tau(\alpha) = \|F\alpha - y\|^2 + \frac{\tau}{2}\|\alpha\|^2,$$

где  $\tau$  — неотрицательный *параметр регуляризации*.

Модифицированное МНК-решение ( $\tau I_n$  — «гребень»):

$$\alpha_\tau^* = (F^T F + \tau I_n)^{-1} F^T y.$$

**Преимущество** сингулярного разложения:

можно подбирать параметр  $\tau$ , вычислив SVD только один раз.

Вектор регуляризованного МНК-решения  $\alpha_\tau^*$   
и МНК-аппроксимация целевого вектора  $F\alpha_\tau^*$ :

$$\alpha_\tau^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y);$$

$$F\alpha_\tau^* = V D U^T \alpha_\tau^* = V \operatorname{diag}\left(\frac{\lambda_j}{\lambda_j + \tau}\right) V^T y = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} v_j (v_j^T y);$$

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

$F\alpha_\tau^* \neq F\alpha^*$ , но зато решение становится гораздо устойчивее.



## Выбор параметра регуляризации $\tau$

Контрольная выборка:  $X^k = (x'_i, y'_i)_{i=1}^k$ ;

$$F'_{k \times n} = \begin{pmatrix} f_1(x'_1) & \dots & f_n(x'_1) \\ \dots & \dots & \dots \\ f_1(x'_k) & \dots & f_n(x'_k) \end{pmatrix}, \quad y'_{k \times 1} = \begin{pmatrix} y'_1 \\ \dots \\ y'_k \end{pmatrix}.$$

Вычисление функционала  $Q$  на контрольных данных  $T$  раз потребует  $O(kn^2 + knT)$  операций:

$$Q(\tau) = \|F' \alpha_\tau^* - y'\|^2 = \left\| \underbrace{F' U}_{k \times n} \operatorname{diag} \left( \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} \right) \underbrace{V^T y}_{n \times 1} - y' \right\|^2.$$

Зависимость  $Q(\tau)$  обычно имеет характерный минимум.

## Регуляризация сокращает «эффективную размерность»

*Сжатие (shrinkage) или сокращение весов (weight decay):*

$$\|\alpha_{\tau}^*\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^T y)^2 < \|\alpha^*\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^T y)^2.$$

Почему говорят о *сокращении эффективной размерности*?

Роль размерности играет след проекционной матрицы:

$$\text{tr } F(F^T F)^{-1} F^T = \text{tr}(F^T F)^{-1} F^T F = \text{tr } I_n = n.$$

При использовании регуляризации:

$$\text{tr } F(F^T F + \tau I_n)^{-1} F^T = \text{tr } \text{diag} \left( \frac{\lambda_j}{\lambda_j + \tau} \right) = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} < n.$$

LASSO — Least Absolute Shrinkage and Selection Operator,  
два эквивалентных варианта постановки задачи:

$$Q(\alpha) = \|F\alpha - y\|^2 \rightarrow \min_{\alpha} \quad \text{при} \quad \sum_{j=1}^n |\alpha_j| \leq \kappa;$$

$$Q(\alpha) = \|F\alpha - y\|^2 + \tau \sum_{j=1}^n |\alpha_j| \rightarrow \min_{\alpha};$$

После замены переменных

$$\begin{cases} \alpha_j = \alpha_j^+ - \alpha_j^-; \\ |\alpha_j| = \alpha_j^+ + \alpha_j^-; \end{cases} \quad \alpha_j^+ \geq 0; \quad \alpha_j^- \geq 0.$$

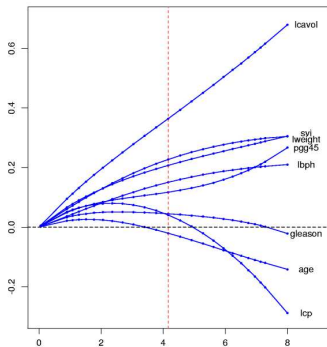
ограничения принимают канонический вид:

$$\sum_{j=1}^n \alpha_j^+ + \alpha_j^- \leq \kappa; \quad \alpha_j^+ \geq 0; \quad \alpha_j^- \geq 0.$$

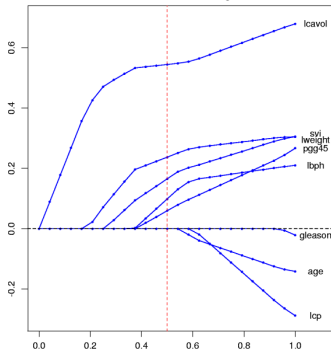
Чем меньше  $\kappa$ , тем больше  $j$  таких, что  $\alpha_j^+ = \alpha_j^- = 0$ .

# Сравнение гребневой регрессии и Лассо

Зависимость  $\{\alpha_j\}$  от  $\frac{1}{\tau}$



Зависимость  $\{\alpha_j\}$  от  $\kappa$



Задача диагностики рака (prostate cancer, UCI)

---

*T.Hastie, R.Tibshirani, J.Friedman. The Elements of Statistical Learning. Springer, 2001.*

- Гребневая регрессия удобно вводится и интерпретируется через сингулярное разложение
- Гребневая регрессия сокращает веса признаков
- LASSO обнуляет веса признаков
- Оба метода имеют параметр регуляризации (селективности), позволяющий определять число признаков (сложность модели) по внешним критериям (по кросс-валидации).

## Метод главных компонент: постановка задачи

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

$g_1(x), \dots, g_m(x)$  — новые числовые признаки,  $m \leq n$ ;

**Требование:** старые признаки должны линейно восстанавливаться по новым:

$$\hat{f}_j(x) = \sum_{s=1}^m g_s(x) u_{js}, \quad j = 1, \dots, n, \quad \forall x \in X,$$

как можно точнее на обучающей выборке  $x_1, \dots, x_\ell$ :

$$\sum_{i=1}^{\ell} \sum_{j=1}^n (\hat{f}_j(x_i) - f_j(x_i))^2 \rightarrow \min_{\{g_s(x_i)\}, \{u_{js}\}}$$

Матрицы «объекты–признаки», старая и новая:

$$F_{\ell \times n} = \begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}; \quad G_{\ell \times m} = \begin{pmatrix} g_1(x_1) & \dots & g_m(x_1) \\ \dots & \dots & \dots \\ g_1(x_\ell) & \dots & g_m(x_\ell) \end{pmatrix}.$$

Матрица линейного преобразования новых признаков в старые:

$$U_{n \times m} = \begin{pmatrix} u_{11} & \dots & u_{1m} \\ \dots & \dots & \dots \\ u_{n1} & \dots & u_{nm} \end{pmatrix}; \quad \hat{F} = GU^T \overset{\text{ХОТИМ}}{\approx} F.$$

**Найти:** и новые признаки  $G$ , и преобразование  $U$ :

$$\sum_{i=1}^{\ell} \sum_{j=1}^n (\hat{f}_j(x_i) - f_j(x_i))^2 = \|GU^T - F\|^2 \rightarrow \min_{G,U},$$

## Теорема

Если  $m \leq \text{rk } F$ , то минимум  $\|GU^T - F\|^2$  достигается, когда столбцы  $U$  — это с.в. матрицы  $F^T F$ , соответствующие  $m$  максимальным с.з.  $\lambda_1, \dots, \lambda_m$ , а матрица  $G = FU$ .

При этом:

- ❶ матрица  $U$  ортонормирована:  $U^T U = I_m$ ;
- ❷ матрица  $G$  ортогональна:  $G^T G = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ ;
- ❸  $U\Lambda = F^T F U$ ;  $G\Lambda = FF^T G$ ;
- ❹  $\|GU^T - F\|^2 = \|F\|^2 - \text{tr } \Lambda = \sum_{j=m+1}^n \lambda_j$ .



Если взять  $m = n$ , то:

❶  $\|GU^T - F\|^2 = 0;$

❷ представление  $\hat{F} = GU^T = F$  точное и совпадает с сингулярным разложением при  $G = V\sqrt{\Lambda}$ :

$$F = GU^T = V\sqrt{\Lambda}U^T; \quad U^TU = I_m; \quad V^TV = I_m.$$

❸ линейное преобразование  $U$  работает в обе стороны:

$$F = GU^T; \quad G = FU.$$

Поскольку новые признаки некоррелированы ( $G^TG = \Lambda$ ), преобразование  $U$  называется *декоррелирующим* (или преобразованием Карунена–Лоэва).

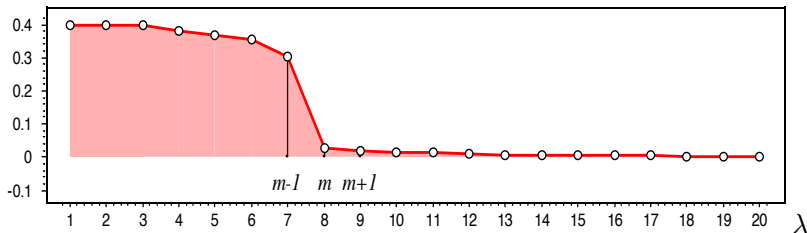
# Эффективная размерность выборки

Упорядочим с.з.  $F^T F$  по убыванию:  $\lambda_1 \geq \dots \geq \lambda_n \geq 0$ .

Эффективная размерность выборки — это наименьшее целое  $m$ , при котором

$$E_m = \frac{\|GU^T - F\|^2}{\|F\|^2} = \frac{\lambda_{m+1} + \dots + \lambda_n}{\lambda_1 + \dots + \lambda_n} \leq \varepsilon.$$

Критерий «крутого склона»: находим  $m$ :  $E_{m-1} \gg E_m$ :



## Решение задачи НК для МЛР в новых признаках

Задача наименьших квадратов для МЛР:  $\|F\alpha - y\|^2 \rightarrow \min_{\alpha}$ .

Заменим  $F$  на её приближение  $G \cdot U^T$ , предполагая  $m \leq n$ :

$$\|G \underbrace{U^T \alpha}_{\beta} - y\|^2 = \|G\beta - y\|^2 \rightarrow \min_{\beta}.$$

Связь нового и старого вектора коэффициентов:

$$\beta = U^T \alpha; \quad \alpha = U\beta.$$

Решение задачи наименьших квадратов относительно  $\beta$  (единственное отличие —  $m$  слагаемых вместо  $n$ ):

$$\beta^* = D^{-1} V^T y; \quad \alpha^* = U D^{-1} V^T y = \sum_{j=1}^m \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y);$$

$$G\beta^* = V V^T y = \sum_{j=1}^m v_j (v_j^T y);$$

- Метод главных компонент позволяет приближать матрицу её низкоранговым разложением.
- Для этого достаточно взять из SVD-разложения первые  $m$  сингулярных чисел и векторов матрицы.
- Этот приём широко используется в анализе данных — в задачах регрессии, классификации, сжатия данных, обработки изображений.