

# Методы машинного обучения

## Многомерная линейная регрессия

Воронцов Константин Вячеславович

[www.MachineLearning.ru/wiki?title=User:Vokov](http://www.MachineLearning.ru/wiki?title=User:Vokov)

вопросы к лектору: [k.vorontsov@iai.msu.ru](mailto:k.vorontsov@iai.msu.ru)

материалы курса:

[github.com/MSU-ML-COURSE/ML-COURSE-24-25](https://github.com/MSU-ML-COURSE/ML-COURSE-24-25)

орг.вопросы по курсу: [ml.cmc@mail.ru](mailto:ml.cmc@mail.ru)

## 1 Многомерная линейная регрессия

- Метод наименьших квадратов
- Многомерная линейная регрессия
- Сингулярное разложение

## 2 Мультиколлинеарность и переобучение

- Проблема мультиколлинеарности
- Число обусловленности матрицы
- Стратегии устранения мультиколлинеарности

## 3 Регуляризация

- $L_2$ -регуляризация: гребневая регрессия
- $L_1$ -регуляризация: лассо Тибширани
- Регуляризаторы для отбора признаков

## Метод наименьших квадратов (МНК)

- $X$  — объекты (часто  $\mathbb{R}^n$ );  $Y$  — ответы (часто  $\mathbb{R}$ , реже  $\mathbb{R}^m$ );  
 $X^\ell = (x_i, y_i)_{i=1}^\ell$  — обучающая выборка;  
 $y_i = y(x_i)$ ,  $y: X \rightarrow Y$  — неизвестная зависимость;
- $a(x, w)$  — модель зависимости,  
 $w \in \mathbb{R}^p$  — вектор параметров модели.
- Метод наименьших квадратов (МНК):

$$Q(w, X^\ell) = \sum_{i=1}^{\ell} \gamma_i (a(x_i, w) - y_i)^2 \rightarrow \min_w,$$

где  $\gamma_i$  — вес, степень важности  $i$ -го объекта.

$Q(w^*, X^\ell)$  — *остаточная сумма квадратов*  
(residual sum of squares, RSS).

# Многомерная линейная регрессия

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

Модель многомерной линейной регрессии:

$$a(x, w) = \sum_{j=1}^n w_j f_j(x), \quad w \in \mathbb{R}^n.$$

Матричные обозначения:

$$F_{\ell \times n} = \begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}, \quad y_{\ell \times 1} = \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}, \quad w_{n \times 1} = \begin{pmatrix} w_1 \\ \dots \\ w_n \end{pmatrix}.$$

Функционал квадрата ошибки:

$$Q(w, X^\ell) = \sum_{i=1}^{\ell} (a(x_i, w) - y_i)^2 = \|Fw - y\|^2 \rightarrow \min_w.$$

## Нормальная система уравнений

Необходимое условие минимума в матричном виде:

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

откуда следует *нормальная система* задачи МНК:

$$F^T F w = F^T y,$$

где  $F^T F$  — матрица размера  $n \times n$ .

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

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

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

## Геометрическая интерпретация МНК

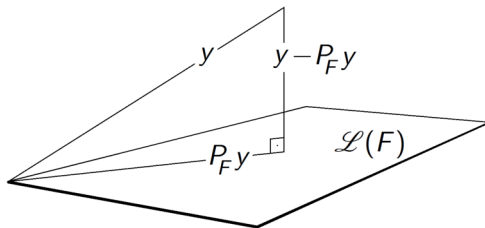
Линейная оболочка столбцов матрицы  $F = (f_1, \dots, f_n)$ ,  $f_j \in \mathbb{R}^\ell$ :

$$\mathcal{L}(F) = \left\{ \sum_{j=1}^n w_j f_j \mid w \in \mathbb{R}^n \right\}$$

$P_F = F(F^T F)^{-1} F^T$  — проекционная матрица

$P_F y$  — проекция вектора  $y \in \mathbb{R}^\ell$  на подпространство  $\mathcal{L}(F)$

$(I_\ell - P_F)y$  — проекция  $y$  на его ортогональное дополнение



МНК — это опускание перпендикуляра в  $\mathbb{R}^\ell$  из  $y$  на  $\mathcal{L}(F)$

## Сингулярное разложение

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  $\ell \times \ell$ -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  $n \times n$ -матрицы  $F^T F$ ;
- ❸  $n \times n$ -матрица  $D$  диагональна,  $D = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$ ,  $\lambda_j \geq 0$  — общие собственные значения матриц  $F^T F$  и  $FF^T$ .

## Решение МНК через сингулярное разложение

Псевдообратная  $F^+ = (F^T F)^{-1} F^T$ , вектор МНК-решения  $w^*$ ,  
МНК-аппроксимация целевого вектора  $Fw^*$ :

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

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

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

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

---

Тождества:  $(AB)^{-1} = B^{-1}A^{-1}$ ,  $(AB)^T = B^T A^T$ ,  $\|w\|^2 = w^T w$



# Мультиколлинеарность и переобучение в линейных моделях

## Возможные причины переобучения:

- слишком мало объектов, слишком много признаков
- линейная зависимость (мультиколлинеарность) признаков:  
линейная модель:  $a(x, w) = \langle w, x \rangle$   
мультиколлинеарность:  $\exists u \in \mathbb{R}^n: \forall x \in X \quad \langle u, x \rangle = 0$   
неединственность решения:  $\forall \gamma \in \mathbb{R} \quad a(x, w) = \langle w + \gamma u, x \rangle$

## Проявления переобучения:

- слишком большие веса  $|w_j|$  разных знаков
- неустойчивость линейной модели  $\langle w, x \rangle$
- $Q(X^\ell) \ll Q(X^k)$

## Простой способ уменьшить переобучение:

- регуляризация  $\|w\| \rightarrow \min$  (сокращение весов, weight decay)

## Неустойчивость модели и число обусловленности матрицы

Если  $\exists \gamma \in \mathbb{R}^n$ :  $\underset{n \times n}{S} \gamma \approx 0$ , то некоторые с.з.  $S$  близки к нулю

*Число обусловленности  $n \times n$ -матрицы  $S$ :*

$$\mu(S) = \|S\| \|S^{-1}\| = \frac{\max_{u: \|u\|=1} \|Su\|}{\min_{u: \|u\|=1} \|Su\|} = \frac{\lambda_{\max}}{\lambda_{\min}}$$

При умножении обратной матрицы на вектор,  $z = S^{-1}u$ , относительная погрешность усиливается в  $\mu(S)$  раз:

$$\frac{\|\delta z\|}{\|z\|} \leq \mu(S) \frac{\|\delta u\|}{\|u\|}$$

В нашем случае:  $w^* = (F^T F)^{-1} F^T y$ ,  $S = F^T F$ ,  $u = F^T y$ , погрешности измерения признаков  $f_j(x_i)$  и ответов  $y_i$  усиливаются в  $\mu(F^T F)$  раз!

## Стратегии устранения мультиколлинеарности

Если матрица  $S = F^T F$  плохо обусловлена, то:

- решение  $w^*$  неустойчиво и плохо интерпретируемо, содержит большие по модулю  $w_j^*$  разных знаков;
- $\|w^*\|$  велико;
- возникает переобучение:  
на обучении  $Q(w^*, X^\ell) = \|Fw^* - y\|^2$  мало;  
на контроле  $Q(w^*, X^k) = \|F'w^* - y'\|^2$  велико;

Стратегии устранения мультиколлинеарности и переобучения:

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

## Гребневая регрессия (ridge regression)

Штраф за увеличение  $L_2$ -нормы вектора весов  $\|w\|$ :

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

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

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

$$\begin{aligned}\frac{\partial Q_\tau(w)}{\partial w} &= 2F^\top(Fw - y) + 2\tau w = 0 \\ w_\tau^* &= (F^\top F + \tau I_n)^{-1} F^\top y.\end{aligned}$$

**Преимущество** сингулярного разложения:

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

## Регуляризованный МНК через сингулярное разложение

Вектор регуляризованного МНК-решения  $w_\tau^*$   
и МНК-аппроксимация целевого вектора  $Fw_\tau^*$ :

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

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

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

$Fw_\tau^* \neq Fw^*$ , но зато решение становится гораздо устойчивее.

## Выбор параметра регуляризации $\tau$

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

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

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

$$Q(w_\tau^*, X^k) = \|F' w_\tau^* - y'\|^2 = \left\| \underbrace{F' U}_{k \times n} \operatorname{diag}\left(\frac{\sqrt{\lambda_j}}{\lambda_j + \tau}\right) \underbrace{V^\top y}_{n \times 1} - y' \right\|^2.$$

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

## Регуляризация сокращает «эффективную размерность»

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

$$\|w_\tau^*\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^\top y)^2 < \|w^*\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^\top y)^2.$$

Почему говорят о *сокращении эффективной размерности*?

Роль размерности играет след проекционной матрицы:

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

При использовании регуляризации:

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

---

Тождество:  $\text{tr}(AB) = \text{tr}(BA)$

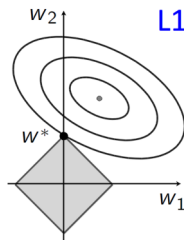
## Регуляризация по $L_1$ -норме для отбора признаков

LASSO — Least Absolute Shrinkage and Selection Operator

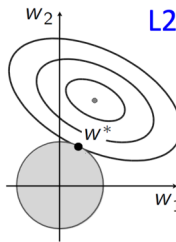
$$\|Fw - y\|^2 + \mu \sum_{j=1}^n |w_j| \rightarrow \min_w \iff \begin{cases} \|Fw - y\|^2 \rightarrow \min_w; \\ \sum_{j=1}^n |w_j| \leq \kappa; \end{cases}$$

LASSO ( $L_1$ ):

$$\sum_{j=1}^n |w_j| \leq \kappa$$



$L_1$



$L_2$

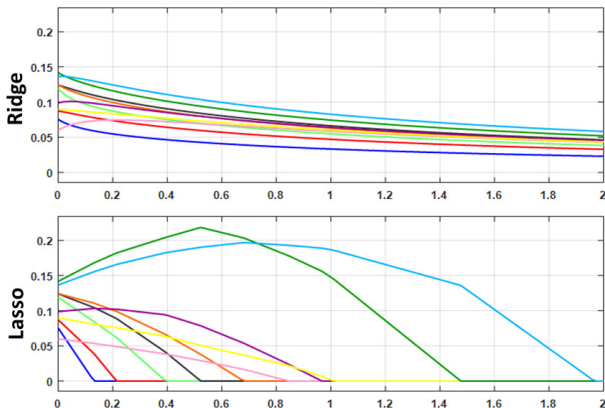
Ridge ( $L_2$ ):

$$\sum_{j=1}^n w_j^2 \leq \kappa$$



## Сравнение $L_2$ (Ridge) и $L_1$ (LASSO) регуляризации

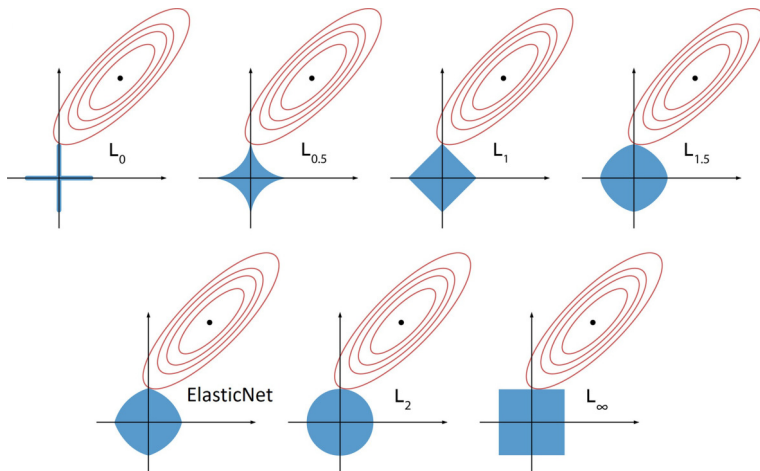
Типичный вид зависимости весов  $w_j$  от селективности  $\mu$



В LASSO с увеличением  $\mu$  усиливается отбор признаков

## Геометрическая интерпретация отбора признаков

Сравнение регуляризаторов по различным  $L_p$ -нормам:



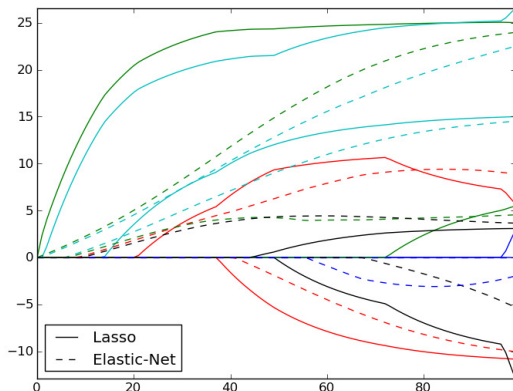
## Двойной регуляризатор $L_1 + L_2$ (Elastic Net)

$$\frac{1}{2} \|Fw - y\|^2 + \mu \sum_{j=1}^n |w_j| + \frac{\tau}{2} \sum_{j=1}^n w_j^2 \rightarrow \min_w$$

- ⊕ Параметр селективности  $\mu$  управляет отбором признаков: чем больше  $\mu$ , тем меньше признаков останется
- ⊕ Есть эффект группировки (grouping effect): значимые зависимые признаки отбираются вместе и имеют примерно равные веса  $w_j$
- ⊖ Приходится подбирать два параметра регуляризации  $\mu$ ,  $\tau$  (есть специальные методы — regularization path)
- ⊖ Шумовые признаки также группируются вместе; по мере увеличения  $\mu$  группы значимых признаков могут отбрасываться, когда ещё не все шумовые отброшены

## Двойной регуляризатор $L_1 + L_2$ (Elastic Net)

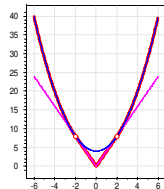
Elastic Net менее жёстко отбирает признаки, чем LASSO.  
Зависимости весов  $w_j$  от коэффициента  $\log \frac{1}{\mu}$ :



## Двойной регуляризатор: склейка $L_1$ и $L_2$

$$\frac{1}{2} \|Fw - y\|^2 + \sum_{j=1}^n R_{\mu}(w_j) \rightarrow \min_w$$

$$R_{\mu}(w_j) = \begin{cases} 2\mu|w_j|, & |w_j| \leq \mu \\ \mu^2 + w_j^2, & |w_j| \geq \mu \end{cases}$$



- ⊕ Только один параметр регуляризации  $\mu$
- ⊕ Отбор признаков с *параметром селективности*  $\mu$
- ⊕ *Эффект группировки*: значимые зависимые признаки ( $|w_j| > \mu$ ) входят в решение совместно (как в Elastic Net)
- ⊕ Шумовые признаки ( $|w_j| < \mu$ ) не группируются и подавляются независимо друг от друга (как в LASSO)

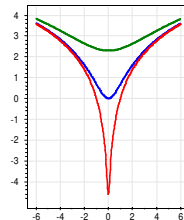
---

Tatarchuk A., Urlov E., Mottl V., Windridge D. A support kernel machine for supervised selective combining of diverse pattern-recognition modalities. 2010.

## Невыпуклый регуляризатор с эффектом отбора признаков

$$\frac{1}{2} \|Fw - y\|^2 + \sum_{j=1}^n \ln(w_j^2 + \frac{1}{\mu}) \rightarrow \min_w$$

$$R(w) = \ln(w^2 + \frac{1}{\mu}) \text{ при } \mu = 0.1, 1, 100$$



- ⊕ Только один параметр регуляризации  $\mu$
- ⊕ Отбор признаков с параметром *селективности*  $\mu$
- ⊕ Есть эффект группировки
- ⊕ Лучше отбирает набор значимых признаков, когда они только совместно обеспечивают хорошее решение

---

Tatarchuk A., Mottl V., Elisseyev A., Windridge D. Selectivity supervision in combining pattern recognition modalities by feature- and kernel-selective Support Vector Machines. 2008.

- Многомерная линейная регрессия
  - через *сингулярное разложение*
- Три приёма против мультиколлинеарности и переобучения:
  - регуляризация
  - отбор признаков
  - преобразование признаков
- $L_2$ -регуляризация, она же гребневая регрессия
  - тоже через *сингулярное разложение*
- $L_1$ -регуляризация (LASSO) и др. негладкие регуляризаторы
  - регулируемый отбор признаков
- Преобразование признаков: *метод главных компонент* и другие *матричные разложения* — в следующем семестре