

Регрессия, регуляризация, работа с признаками

Михайлов Дмитрий, Смирнов Иван

Санкт-Петербургский государственный университет
Кафедра статистического моделирования

Санкт-Петербург, 2023

Идея

Имеется набор данных (обучающая выборка)

$$\mathbf{X} \in \mathbb{R}^{n \times p}, \quad \mathbf{y} \in \mathbb{R}^n$$

$\mathbf{x}_i \in \mathbb{R}^p$ — вектор-строки \mathbf{X} , $X_i \in \mathbb{R}^n$ — вектор-столбцы \mathbf{X} .

Задача: уметь предсказывать y (ответ) по новым \mathbf{x}_i (объектам), установив некоторую зависимость на обучающей выборке.

Постановка задачи

"До эксперимента":

$\xi \in \mathbb{R}^p$ — случайный вектор, $\eta, \varepsilon \in \mathbb{R}$ — случайные величины.
Предполагаем, что η и ξ функционально зависимы:

$$\eta = \varphi(\xi) + \varepsilon$$

Обычно $E\varepsilon = 0$, $D\varepsilon = \sigma^2$, $\xi \perp \varepsilon$.

"После эксперимента":

$\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathcal{L}(\xi)$; $y_1, \dots, y_n \sim \mathcal{L}(\eta)$ — выборки, которые наблюдаем. Модель для всех $i \in 1 : n$

$$y_i = \varphi(\mathbf{x}_i) + \varepsilon_i$$

Задача: найти функцию φ .

Этапы обучения модели

Модель

$$y_i = \varphi(\mathbf{x}_i) + \varepsilon_i$$

Задача: найти функцию φ .

1. Выбор модели регрессии (класс рассматриваемых $\varphi(\cdot)$)

Линейная модель: $\varphi(\mathbf{x}_i, \beta) = \sum_{j=1}^p \beta_j \mathbf{x}_i[j]$, $i \in 1 : n$

2. Выбор функции потерь (loss function)

Квадратичная функция потерь: $\sum_{i=1}^n (y_i - \varphi(\mathbf{x}_i, \beta))^2$

3. Выбор метода обучения (training)

МНК: $\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \varphi(\mathbf{x}_i, \beta))^2$

4. Выбор метода проверки (test)

MSE: $\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} (y_i^{\text{test}} - \varphi(\mathbf{x}_i^{\text{test}}, \hat{\beta}))^2$

Оптимизация. Матричный вид в общем случае

- ▶ $\mathbf{X} \in \mathbb{R}^{n \times p}$ — матрица данных
- ▶ $\mathbf{y} \in \mathbb{R}^n$ — вектор ответов
- ▶ $\boldsymbol{\beta} \in \mathbb{R}^d$ — вектор параметров
- ▶ $\boldsymbol{\varphi}(\mathbf{X}, \boldsymbol{\beta}) := (\varphi(\mathbf{x}_1, \boldsymbol{\beta}), \dots, \varphi(\mathbf{x}_n, \boldsymbol{\beta}))^T$ — функция от выборки и параметров
- ▶ $\mathcal{L}(\boldsymbol{\varphi}(\mathbf{X}, \boldsymbol{\beta}), \mathbf{y})$ — некоторая функция потерь

Решение задачи регрессии — вектор коэффициентов $\hat{\boldsymbol{\beta}}$.

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \mathcal{L}(\boldsymbol{\varphi}(\mathbf{X}, \boldsymbol{\beta}), \mathbf{y})$$

Выбор функции потерь

Какие варианты есть?

- ▶ $\|\varphi(\mathbf{X}, \beta) - \mathbf{y}\|_2^2$ — квадратичная ошибка (l_2 -норма)
- ▶ $\|\varphi(\mathbf{X}, \beta) - \mathbf{y}\|_1$ — модуль ошибки (l_1 -норма)

Как выбирать?

- ▶ Предположения о распределении остатков ϵ
- ▶ Простота функции \mathcal{L} для оптимизации
- ▶ Точность данных/наличие выбросов

Линейная регрессия: постановка

Модель: $y = X\beta + \varepsilon$

- ▶ $y \in \mathbb{R}^n$ — вектор ответов, $\varepsilon \in \mathbb{R}^n$ — вектор ошибок, $E\varepsilon = 0$
- ▶ $X \in \mathbb{R}^{n \times p}$ — матрица данных
- ▶ $\beta \in \mathbb{R}^p$ — вектор параметров

Решение задачи линейной регрессии — вектор $\hat{\beta}$.

Классическая задача (с квадратичной функцией потерь):

$$\hat{\beta} = \arg \min_{\beta} \|y - X\beta\|_2^2$$

Ищем наилучшую аппроксимацию в подпространстве столбцов X .

МНК-оценка и её особенности I

Задача

$$\hat{\beta} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

Решение

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- ✓ Явный вид решения
- ✓ Простота функции \mathcal{L} для оптимизации
 - ! Точность данных/наличие выбросов
 - ! Конкретные предположения о распределении остатков ϵ
 - ! Мультиколлинеарность
 - ! $n \geq p$, иначе решений бесконечно много

МНК-оценка и её особенности II

Для оценок $\hat{\beta}$ имеет место разложение

$$\text{MSE} = E(\beta - \hat{\beta})^2 = \underbrace{D\hat{\beta}}_{\text{дисперсия}} + \underbrace{(E\hat{\beta} - \beta)^2}_{\text{смещение}}$$

Рассмотрим МНК-оценку $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

- ▶ $E\hat{\beta} = \beta$, то есть оценка несмещённая
- ▶ $D\hat{\beta} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ (в случае, когда $\varepsilon_i \sim N(0, \sigma^2)$)
- ▶ В классе несмещённых оценок $\hat{\beta}$ обладает наименьшей дисперсией ($\hat{\beta}$ — BLUE)
- ▶ Если остатки распределены нормально, то $\hat{\beta}$ — ОМП

Анализ остатков

Модель: $y = X\beta + \varepsilon$

- ▶ Распределение ε — нормальное, $E\varepsilon = \mathbf{0}$:
 - ▶ Классическая оценка по МНК — ОМП и BLUE
 - ▶ Работают стандартные критерии значимости
- ▶ Распределение ε неизвестно или с тяжёлыми хвостами:
 - ▶ Оценка по МНК уже не ОМП и не BLUE

МНК-оценка: вычислительный взгляд

Решение

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- ▶ Необходимо вычислить $\hat{\beta}$
- ▶ Сингулярное разложение: $\mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{U}^T$
 - ▶ \mathbf{V} и \mathbf{U} — ортогональные, \mathbf{D} — диагональная
 - ▶ $\mathbf{V} = (V_1, V_2, \dots, V_n) \in \mathbb{R}^{n \times n}$, V_i — с. векторы $\mathbf{X} \mathbf{X}^T$
 - ▶ $\mathbf{U} = (U_1, U_2, \dots, U_p) \in \mathbb{R}^{p \times p}$, U_i — с. векторы $\mathbf{X}^T \mathbf{X}$
 - ▶ $\mathbf{D} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$, $\lambda_j \geq 0$ — с. значения $\mathbf{X}^T \mathbf{X}$

Отсюда $\hat{\beta} = \mathbf{U} \mathbf{D}^{-1} \mathbf{V}^T \mathbf{y}$

- ▶ $\mathbf{D}^{-1} = \text{diag}(1/\sqrt{\lambda_1}, \dots, 1/\sqrt{\lambda_n})$

Борьба с мультиколлинеарностью

Мультиколлинеарность ($X^T X$ плохо обусловлена) влечёт

- ▶ Неустойчивость решения
- ▶ Высокая дисперсия у $\hat{\beta} \Rightarrow$ высокая MSE

Возможные решения проблемы мультиколлинеарности

- ▶ Уменьшение числа признаков (отбор признаков)
- ▶ Регуляризация
- ▶ Преобразование признаков (АГК)

Регуляризация. Гребневая регрессия

Модель: $y = X\beta + \epsilon$

Вернёмся к разложению MSE

$$\text{MSE} = E(\beta - \hat{\beta})^2 = \underbrace{D\hat{\beta}}_{\text{дисперсия}} + \underbrace{(E\hat{\beta} - \beta)^2}_{\text{смещение}}$$

Несмещённая оценка может иметь большую дисперсию и MSE.
Возьмём оценку по МНК и сделаем её смещённой:

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y, \quad \lambda > 0$$

Используем SVD и получим

$$\hat{\beta}_{\text{ridge}} = \sum_{i=1}^n \frac{\lambda_j}{\lambda_j + \lambda} U_j (V_j^T y)$$

Отделили знаменатель от нуля. Устойчивость вычислений повышается.

Гребневая регрессия. Дисперсия

$$\text{MSE} = E(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^2 = \underbrace{D\hat{\boldsymbol{\beta}}}_{\text{дисперсия}} + \underbrace{(E\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^2}_{\text{смещение}}$$

Оценка по методу гребневой регрессии

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}, \quad \lambda > 0$$

Смещение контролируется параметром λ .

Дисперсия:

- ▶ $\hat{\boldsymbol{\beta}}_{\text{ridge}} = \sum_{i=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \lambda} U_j (V_j^T \mathbf{y})$
- ▶ λ_j убывают
- ▶ $\frac{\sqrt{\lambda_j}}{\lambda_j + \lambda}$ штрафуют компоненты с наименьшей дисперсией
- ▶ $D\hat{\boldsymbol{\beta}}$ уменьшается $\Rightarrow \text{MSE} \downarrow$

Гребневая регрессия как задача оптимизации

Есть две эквивалентные формулировки:

Ridge Regression

$$\hat{\beta}_{\text{ridge}}(\lambda_1) = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_2^2$$

$$\hat{\beta}_{\text{ridge}}(\lambda_2) = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2, \text{ s.t. } \|\beta\|_2^2 \leq \lambda_2$$

Регуляризация. Lasso

Есть две эквивалентные формулировки (явного вида нет):

Lasso

$$\hat{\beta}_{\text{Lasso}}(\lambda_1) = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1^2$$

$$\hat{\beta}_{\text{Lasso}}(\lambda_2) = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2, \text{ s.t. } \|\beta\|_1^2 \leq \lambda_2$$

Особенности:

- ▶ Уменьшение MSE
- ▶ Интерпретируемость результатов
- ▶ Выбор параметра: кросс-валидация

Ridge VS Lasso

Тут мы видим, что у lasso точка пересечения $(0; \beta_2^{lasso})$, а у ridge $(\beta_1^{ridge}; \beta_2^{ridge})$, соответственно lasso позволяет нам таким образом упростить модель и повысить её интерпретируемость.

Dimension Reduction of Feature Space with LASSO

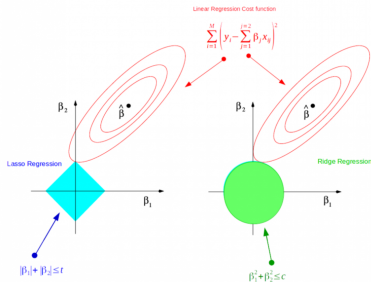
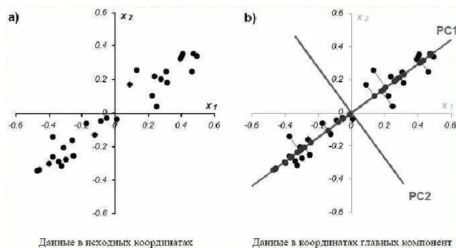


Figure 3: Why LASSO can reduce dimension of feature space? Example on 2D feature space. Modified from the plot used in 'The Elements of Statistical Learning' by Author.

Преобразование признаков. АГК

В методе главных компонент строится минимальное число новых признаков, по которым исходные признаки восстанавливаются линейным преобразованием с минимальными погрешностями.



Отбор признаков. Общий подход к выбору модели

Предположим, что есть некоторое семейство построенных моделей $\{\mathcal{M}_i\}_{i \in I}$.

Хотим выбрать лучшую модель \mathcal{M}^* для предсказания.

Классические методы выбора модели

- ▶ Кросс-валидация
- ▶ Информационные критерии
 - ▶ $\text{AIC}_i = 2p_i - 2 \ln \mathcal{L}(\mathbf{X}; \mathcal{M}_i)$
 - ▶ $\text{BIC}_i = p_i \ln n - 2 \ln \mathcal{L}(\mathbf{X}; \mathcal{M}_i)$