

Регрессия, регуляризация, отбор признаков

Михаил Козак, Павел Мехнин, Данил Шкурат

Санкт-Петербургский государственный университет
Кафедра статистического моделирования

Семинар по статистическому и машинному обучению



Санкт-Петербург, 2022

Обучение с учителем — это направление машинного обучения, объединяющее алгоритмы и методы построения моделей на основе совокупности прецедентов (обучающей выборки), содержащих пары «известный вход — известный выход».

- X — множество объектов, заданное своими признаками (точки p -мерного пространства)
- Y — множество ответов (действительные числа)
- Предполагаем наличие неизвестной зависимости между объектами и ответами $y : X \rightarrow Y$
- Обучающая выборка из множества объектов $\{x_1, x_2, \dots, x_n\} \subset X$ и известных ответов (откликов) $y_i = y(x_i) \in Y, i = 1, \dots, n$

Задача:

Найти $y^* : X \rightarrow Y$ — отображение, приближающее неизвестную функцию y на всём множестве X , то есть восстановить зависимость, способную для любого объекта выдать достаточно точный отклик.

Задача регрессии как задача оптимизации

- $X_n = (x_i, y_i)_{i=1}^n$ — обучающая выборка, $x_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}$.
- $y_i = y(x_i) \in Y, i = 1, \dots, n$
- Модель регрессии: параметрическое семейство функций $f(x, \beta)$, где $\beta \in B \subset \mathbb{R}^p$ — вектор параметров модели.
- Средняя квадратичная ошибка (функционал качества, наиболее часто применяющийся в задачах регрессии):

$$Q(\theta, X_n) = \sum_{i=1}^n (f(x_i, \beta) - y_i)^2.$$

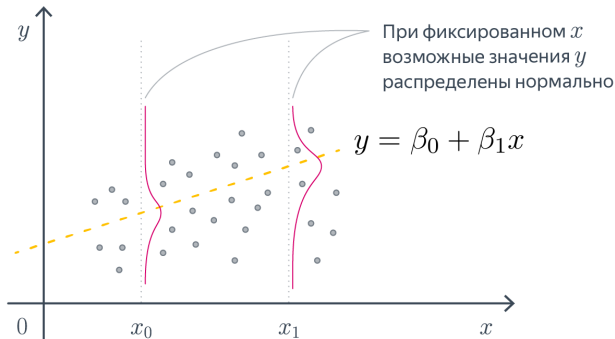
- Задача обучения по МНК — задача минимизации

$$Q(\beta, X_n) \rightarrow \min_{\beta \in B}.$$

Линейная регрессия: постановка

Модель: $y = X\beta + \epsilon$

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



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

Задача оптимизации (с квадратичной функцией потерь):

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

Полученную оценку называют оценкой по методу наименьших квадратов. Она имеет явный вид:

$$\begin{aligned}\hat{\beta}_{\text{МНК}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^+ \mathbf{y} \\ \hat{\mathbf{y}} &= \mathbf{X} \hat{\beta}_{\text{МНК}}\end{aligned}$$

Теорема Гаусса–Маркова утверждает, что $\hat{\beta}_{\text{МНК}}$ имеет наименьшую дисперсию среди всех несмещённых оценок (best linear unbiased estimate — BLUE).

Достаточно быстро вычисляется посредством применения сингулярного разложения.

Сингулярное разложение: $X = VDU^T$

- V и U — ортогональные, D — диагональная
- $V = (V_1, V_2, \dots, V_n) \in \mathbb{R}^{n \times n}$, V_i — с. векторы XX^T
- $U = (U_1, U_2, \dots, U_n) \in \mathbb{R}^{p \times n}$, U_i — с. векторы $X^T X$
- $D = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$, $\lambda_j \geq 0$ — с. значения $X^T X$

Подставим в формулу для $\hat{\beta}_{\text{МНК}}$ вместо матрицы X её сингулярное разложение и получим

$$\hat{\beta} = (X^T X)^{-1} X^T y = (UDV^T VDU^T)^{-1} UDV^T y = UD^{-1} V^T y,$$

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

$$\hat{\beta}_{\text{МНК}} = \sum_{j=1}^p \frac{1}{\sqrt{\lambda_j}} U_j (V_j^T y)$$

Пусть матрица данных содержит несколько сильно коррелированных признаков (есть $\lambda_j \rightarrow 0$).

Что будет происходить в таком случае с МНК-оценкой?

- Решение $\hat{\beta}$ неустойчиво
- Решение неинтерпретируемо, $\|\hat{\beta}\| \rightarrow \infty$
- Высокая дисперсия у $\hat{\beta} \rightarrow$ высокая MSE
- Ответы на контрольной выборке неустойчивы (переобучение)

Способы решения проблемы:

- Регуляризация
- Уменьшение числа признаков (отбор признаков)
- Преобразование признаков (анализ главных компонент)

Хорошая оценка $\hat{\beta}$ должна иметь низкую среднеквадратическую ошибку

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

Несмещенная МНК-оценка не гарантирует минимизацию всей MSE, т.к. может иметь большую дисперсию (если матрица X близка к вырожденной).

Введение небольшого смещения в оценке может привести к уменьшению дисперсии и тем самым уменьшению MSE_{test} .

- Вводим штраф за увеличение нормы вектора β и переходим к минимизации следующей функции:

$$Q_{\tau}(\beta) = ||\mathbf{X}\beta - \mathbf{y}||^2 + \tau||\beta||^2 \rightarrow \min_{\beta},$$

где τ — неотрицательный параметр регуляризации.

- В развернутом виде задача оптимизации записывается так:

$$\sum_{i=1}^n \left(y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \tau \sum_{j=1}^p \beta_j^2 \rightarrow \min_{\beta}.$$

- Решение:

$$\hat{\beta}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \tau \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}.$$

Используем SVD и получим

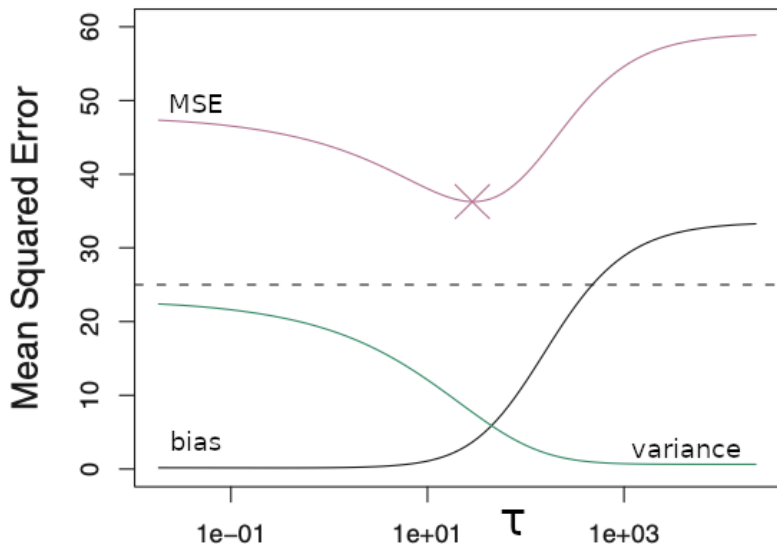
$$\hat{\beta}_{\text{ridge}} = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} U_j (V_j^T \mathbf{y})$$

Отделили знаменатель от нуля. Устойчивость вычислений повышается.

Преимущество SVD разложения: можно подбирать параметр τ вычислив SVD только один раз.

Чем больше коэффициент регуляризации τ , тем устойчивее решение, но больше смещение.

Когда $\tau = 0$, то гребневая регрессия совпадает с обычной регрессией, но при $\tau \rightarrow \infty$ коэффициенты регрессии стремятся к нулю.



Скользящий контроль:

- выбираем сетку значений τ ;
- вычисляем ошибку кросс-проверки для каждого значения τ ;
- выбираем τ с наименьшим значением ошибки кросс-проверки;
- перестраиваем модель со всеми наблюдениями с выбранным значением τ .

Эвристика

Скользящий контроль — вычислительно трудоёмкая процедура. Известна практическая рекомендация брать τ в отрезке $[0.1, 0.4]$, если столбцы матрицы X заранее стандартизованы.

- В качестве штрафа за увеличение нормы вектора β используется его l_1 -норма.
- Метод LASSO решает следующую задачу минимизации:

$$\|X\beta - y\|_2^2 + \tau\|\beta\|_1^2 \rightarrow \min_{\beta},$$

где τ — неотрицательный параметр регуляризации.

- Задача оптимизации в развернутом виде:

$$\sum_{i=1}^n \left(y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \tau \sum_{j=1}^p |\beta_j| \rightarrow \min_{\beta_1, \dots, \beta_p}.$$

Задачу lasso-оптимизации можно переписать в форме с ограничениями:

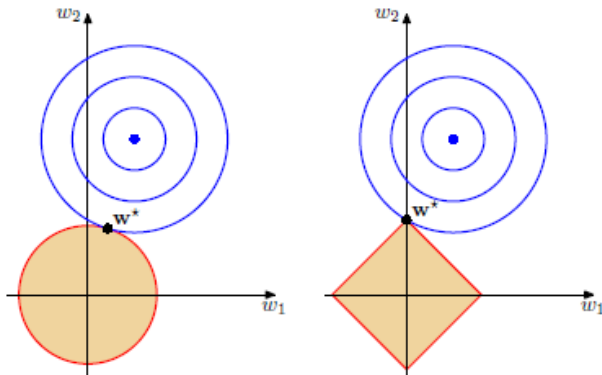
$$\begin{cases} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \rightarrow \min_{\beta_1, \dots, \beta_p}, \\ \sum_{j=1}^p |\beta_j| \leq \varkappa, \end{cases}$$

где $\varkappa = 1/\tau$.

Особенности:

- Уменьшение MSE
- Интерпретируемость результатов
- Быстрое вычисление $\hat{\beta}_{\text{Lasso}}(\tau)$
- Выбор параметра: кросс-валидация

Сравнение гребневой регрессии и Лассо



Синие линии уровня функционала качества (синяя точка — безусловный минимум, который достигается на МНК решении).

Оранжевая зона — ограничения, задаваемые L2 и L1-регуляризаторами. Чёрная точка — минимум целевой функции при заданном ограничении.

- Оба метода успешно решают проблему мультиколлинеарности
- Гребневая регрессия использует все признаки
- Лассо производит отбор признаков, что предпочтительнее, если среди признаков есть шумовые или измерения признаков связаны с ощутимыми затратами.
- С помощью кросс-валидации можно определить какой подход лучше для конкретных данных.

Решается задача оптимизации

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \tau_1 \|\boldsymbol{\beta}\|_1^2 + \tau_2 \|\boldsymbol{\beta}\|_2^2 \rightarrow \min_{\boldsymbol{\beta}}.$$

- ✓ Elastic net — это комбинация методов Lasso и Ridge:
 - Когда $\tau_1 = 0$: Ridge регрессия
 - Когда $\tau_2 = 0$: Lasso регрессия
- ✓ Elastic net в целом лучше, чем Lasso при наличии коррелированных признаков
- ✓ В отличие от Ridge регрессии, когда $p > n$, Elastic net может учитывать более n переменных

Алгоритм полного перебора (Best Subset Selection)

Преимущества:

- ✓ простота реализации
- ✓ гарантированный результат
- ✓ эффективен, когда информативных признаков немного (≤ 5) и общее количество признаков также не велико ($\leq 20 \dots 100$)

Недостатки:

- ✗ в общем случае очень долго: $O(2^p)$. Например, для $p = 20$: $2^p = 1,048,576$ моделей
- ✗ чем больше перебирается вариантов, тем больше перобучение

Решение: эвристические алгоритмы сокращенного перебора

Жадные (greedy) алгоритмы

- Forward Stepwise Selection

Начинаем со свободного члена, потом добавляем на каждом шаге предиктор, который максимально уменьшает ошибку.

Подмножества получаются вложенные — для $p = 20$:
 $p(p+1)/2 = 210$ моделей.

- Backward Stepwise Selection

Начинаем с полной регрессии и на каждом шаге убираем предиктор, который оказывает меньше всего влияния на ошибку.