

# Регуляризация в регрессии

Ширинкина Дарья Андреевна, гр. 622

Санкт-Петербургский государственный университет  
Математико-механический факультет  
Статистическое моделирование

Санкт-Петербург  
2017г.

## Пусть

- $x_1, \dots, x_n \in \mathbb{R}^p$  — независимые одинаково распределенные случайные величины;
- $X = [X_1, \dots, X_p]$ , где  $X_i = (x_{1i}, \dots, x_{ni})^T$ ,  $i = 1, \dots, p$ .

Предполагаем существование неизвестной  $f$  такой, что

$$y_i = f(x_i) + \varepsilon_i,$$

где

- $Y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ ;
- $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$  — независимые случайные величины;
- $\varepsilon_i$  и  $x_j$  независимы для  $\forall i, j$ ;
- $\mathbb{E}\varepsilon_i = 0$ ,  $i = 1, \dots, n$  и  $\mathbb{E}\varepsilon_i^2 = \sigma^2$ .

**Задача:** Оценить функцию  $f$ .

## Обучающая выборка:

- $x_1, \dots, x_n$  — выборка, участвующая в оценке функции  $f$  (обучающая выборка);
- $y_i = f(x_i) + \varepsilon_i, i = 1, \dots, n$ .

## Тестовая выборка:

- $x'_1, \dots, x'_k$  — выборка, по которой оценивается качество оценки функции  $f$  (тестовая выборка);
- $y'_i = f(x'_i) + \varepsilon'_i, i = 1, \dots, k$ .

Считаем, что  $Y = (y_1, \dots, y_n)^T$  и  $X$  — центрированы.

**Модель многомерной линейной регрессии:**

$$y_i = f(x_i, \beta) + \varepsilon_i = \sum_{j=1}^p \beta_j x_{ij} + \varepsilon_i.$$

**Задача минимизации:**

$$\text{MSE}_{\text{training}} = \frac{1}{n} \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 \rightarrow \min_{\beta}.$$

**Решение МНК:**  $\hat{\beta} = (X^T X)^{-1} X^T Y$ .

Насколько хорошо предсказывает  $\hat{f}(x) = \sum_{i=1}^p \hat{\beta}_i x_i$ ?

**Проблема:** минимизируем  $\text{MSE}_{\text{training}}$ , но хотим минимизировать

$$\text{MSE}_{\text{test}} = \frac{1}{n} \sum_{i=1}^n (y_i' - \sum_{j=1}^p \beta_j x_{ij}')^2.$$

- Нет гарантии, что минимум  $\text{MSE}_{\text{training}}$  будет соответствовать минимуму  $\text{MSE}_{\text{test}}$ .
- Когда  $\text{MSE}_{\text{test}} \gg \text{MSE}_{\text{training}}$ , говорят, что происходит переобучение.

Пусть

- $x'_i$  — реализация случайной величина из тестовой выборки;
- $y'_i = f(x'_i) + \varepsilon'_i$  — известное значение.

$$\mathbb{E}(y'_i - \hat{f}(x'_i))^2 = \text{Var}(\hat{f}(x'_i)) + (\text{Bias}(\hat{f}(x'_i)))^2 + \text{Var}(\varepsilon'_i).$$

- Как правило, при увеличении сложности метода (увеличение числа параметров) дисперсия будет увеличиваться, а смещение будет уменьшаться.
- Введение небольшого смещения в оценке может привести к значительному уменьшению дисперсии и тем самым уменьшению  $\text{MSE}_{\text{test}}$ .

**Регуляризация:** вводим ограничения на коэффициенты  $\beta$ .

Для чего используем регуляризацию:

- можем уменьшить дисперсию оценки за счет введения смещения и тем самым уменьшить  $MSE_{\text{test}}$  (особенно, когда  $p > n$ );
- можем производить отбор значимых признаков, делая коэффициенты при них равными нулю.

Задача минимизации:

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

где  $\lambda \geq 0$  — неотрицательный параметр регуляризации (tuning parameter).

- $\lambda \sum_{j=1}^p \beta_j^2$  мало, когда  $\beta_1, \dots, \beta_p$  близки к нулю.
- Когда  $\lambda = 0$ , то гребневая регрессия совпадает с обычной регрессией, но при  $\lambda \rightarrow \infty$  коэффициенты регрессии стремятся к нулю.
- Необходимо выбрать хорошее значение  $\lambda$ .



# Способ решения оптимизационной задачи

Модифицированное МНК решение гребневой регрессии:

$$\hat{\beta}_{\lambda}^R = (X^T X + \lambda I_p)^{-1} X^T Y.$$

Решение через сингулярное разложение, где  $X = VDU^T$ :

- МНК

$$\hat{\beta} = \sum_{j=1}^p \frac{1}{\sqrt{\lambda_j}} U_j (V_j^T Y);$$

- МНК гребневой регрессии

$$\hat{\beta}_{\lambda}^R = U(D^2 + \lambda I_p)^{-1} D V^T y = \sum_{j=1}^p \frac{\sqrt{\lambda_j}}{\lambda_j + \lambda} U_j (V_j^T Y).$$

С помощью сингулярного разложения можно быстро выбирать параметр  $\lambda$ .

Как выбрать параметр  $\lambda$ :

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

- $\beta = (\beta_1, \dots, \beta_p)^T$  имеет априорное распределение  $p(\beta)$ ;
- $f(Y|X, \beta)$  — функция правдоподобия исходных данных.

При фиксированном  $X$  апостериорное распределение  $p(\beta|X, Y)$  пропорционально

$$f(Y|X, \beta)p(\beta|X) = f(Y|X, \beta)p(\beta).$$

Предполагая, что

- 1 линейная модель имеет независимые и нормально распределенные ошибки;
- 2  $p(\beta) = \prod_{j=1}^p g(\beta_j)$  для некоторой плотности  $g$ .

Если  $g$  — плотность  $N(0, \lambda)$ , то оценка апостериорного максимума  $\beta$  совпадает с решением гребневой регрессии.

Пусть  $X^T X = \Sigma$  и  $Y = (y_1, \dots, y_n)^T$ .

Оценка гребневой регрессии через МНК оценку:

$$\hat{\beta}_\lambda^R = (I_p + \lambda \Sigma^{-1}) \hat{\beta}.$$

Оценка гребневой регрессии имеет смещение:

$$\begin{aligned} \mathbb{E} \hat{\beta}_\lambda^R &= \mathbb{E}[(I_p + \lambda \Sigma^{-1}) \hat{\beta}] = \\ &= (I_p + \lambda \Sigma^{-1}) \beta. \end{aligned}$$

Если  $\lambda = 0$ , то оценка гребневой регрессии не имеет смещения.

- Оценки МНК инварианты относительно умножения признака на константу, то есть значение  $f(x_j)\hat{\beta}_j$  не зависит от масштаба  $j$ -го признака.
- Инвариант относительно масштаба теряется в случае гребневой регрессии, оценки МНК гребневой регрессии могут сильно измениться при умножении заданного признака на константу.

**Вывод:** гребневую регрессию нужно использовать после стандартизации признаков.

## Проблема:

- в конечную модель входят все начальные признаки;
- если признаков много, то усложняется интерпретация.

Задача минимизации:

$$\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j| \rightarrow \min_{\beta},$$

где  $\lambda \geq 0$  — неотрицательный параметр регуляризации (tuning parameter).

- Как и в гребневой регрессии  $\lambda \sum_{j=1}^p |\beta_j|$  мало, когда  $\beta_1, \dots, \beta_p$  близки к нулю.
- При увеличении параметра  $\lambda$  некоторые коэффициенты регрессии становятся равными нулю.
- Как и в гребневой регрессии необходимо выбрать хорошее значение  $\lambda$ .

# Способ решения оптимизационной задачи

Задача:

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

эквивалента задаче минимизации с ограничением:

$$\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 \rightarrow \min_{\beta}, \quad \sum_{j=1}^p |\beta_j| \leq s,$$

где параметру  $\lambda$  соответствует параметр  $s$ .

- Чем меньше  $s$ , тем больше нулевых значений коэффициентов  $\beta$ .
- Значение параметра  $\lambda$  выбирается как в гребневой регрессии с помощью кросс-проверки.

# Вероятностная интерпретация

- $\beta = (\beta_1, \dots, \beta_p)^T$  имеет априорное распределение  $p(\beta)$ ;
- $f(Y|X, \beta)$  — функция правдоподобия исходных данных.

При фиксированном  $X$  апостериорное распределение  $p(\beta|X, Y)$  пропорционально

$$f(Y|X, \beta)p(\beta|X) = f(Y|X, \beta)p(\beta).$$

Предполагая, что

- 1 линейная модель имеет независимые и нормально распределенные ошибки;
- 2  $p(\beta) = \prod_{j=1}^p g(\beta_j)$  для некоторой плотности  $g$ .

Если  $g$  — плотность распределения Лапласа с нулевым средним и параметром масштаба  $\lambda$ , то оценка апостериорного максимума  $\beta$  является решением Лассо.



Почему Лассо обнуляет коэффициенты.

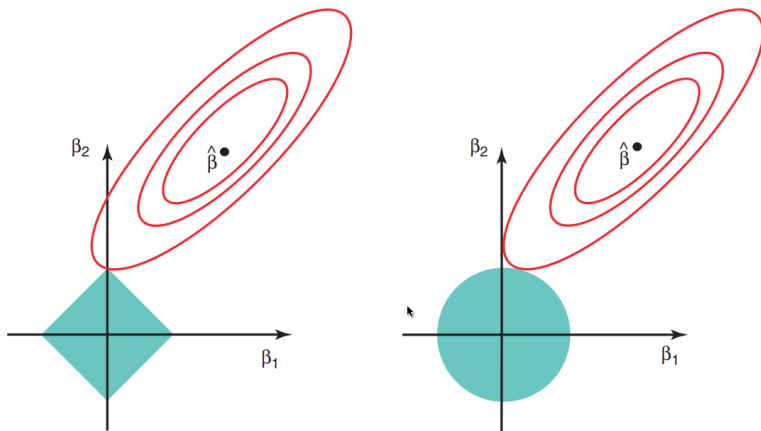


Рис.: Границы ошибки  $\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2$  и ограничений  $\sum_{j=1}^p |\beta_j| \leq s$  для Лассо (слева) и  $\sum_{j=1}^p \beta_j^2 \leq s$  для гребневой регрессии (справа).

# Сравнение гребневой регрессии и Лассо

Рассмотрим простой случай, когда  $n = p$  и  $X$  — диагональная матрица с 1 на диагонали.

**МНК:**

$$\sum_{j=1}^p (y_j - \beta_j)^2.$$

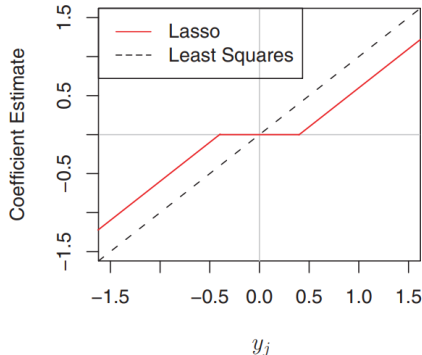
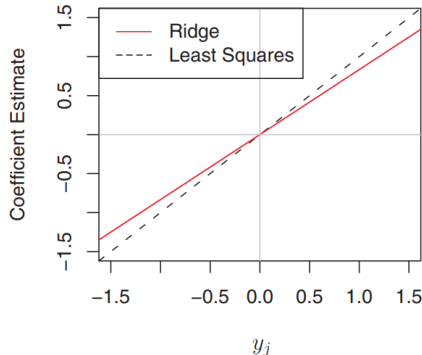
**Решение гребневой регрессии:**

$$\hat{\beta}_{\lambda}^R = \frac{y_j}{1 + \lambda}.$$

**Решение Лассо:**

$$\hat{\beta}_{\lambda}^L = \begin{cases} y_j - \lambda/2, & y_j > \lambda/2 \\ y_j + \lambda/2, & y_j < -\lambda/2 \\ 0, & |y_j| \leq \lambda/2. \end{cases}$$

# Сравнение гребневой регрессии и Лассо



- Гребневая регрессия уменьшает каждый коэффициент с равной пропорцией;
- Лассо уменьшает значения коэффициентов на одинаковое значение;
- В Лассо, если коэффициент по модулю меньше  $\lambda/2$ , то его значение становится равным нулю.

# Сравнение гребневой регрессии и Лассо

Нельзя выделить ни одну из моделей (Лассо или гребневая регрессия) как лучшую.

Можно ожидать, что

- Лассо будет иметь ошибку меньше, когда в модели мало значимых признаков (коэффициенты при таких признаках будут равны нулю);
- Гребневая регрессия будет иметь ошибку меньше, когда  $Y$  будет зависеть от признаков, которые имеют примерно равную значимость.

С помощью кросс-проверки можно определить какой подход лучше.