

Санкт-Петербургский государственный университет
Прикладная математика и информатика
Кафедра статистического моделирования

Ширинкина Дарья Андреевна

Регуляризация в регрессии

Санкт-Петербург
2017

Содержание

1	Задача регрессии	3
1.1	Задача регрессии на генеральном языке	3
1.2	Задача регрессии на выборочном языке	3
1.3	Линейная регрессия и МНК	3
1.4	Проблема МНК	4
2	Гребневая регрессия	5
2.1	Задача гребневой регрессии	5
2.2	Решение задачи минимизации и выбор параметра регуляризации	5
2.3	Свойства и проблемы	6
3	Лассо	7
3.1	Задача минимизации Лассо	7
3.2	Решение задачи минимизации и выбор параметра регуляризации	7
3.3	Свойства	7
4	Сравнение гребневой регрессии и Лассо	9

1 Задача регрессии

1.1 Задача регрессии на генеральном языке

Предполагаем, что существует неизвестная зависимость $f : \mathbb{R}^p \rightarrow \mathbb{R}$ между случайными величинами $\eta \in \mathbb{R}$ и $\xi \in \mathbb{R}^p$:

$$\eta = f(\xi) + \epsilon,$$

где ϵ — случайная величина (ошибка) такая, что ξ и ϵ независимы. Будем считать также, что $\mathbb{E}\epsilon = 0$ и $\mathbb{E}\epsilon^2 = \sigma^2$.

Задача: хотим оценить функцию f .

1.2 Задача регрессии на выборочном языке

На практике имеем выборку размера n из случайной величины ξ : $X = [X_1, \dots, X_p]$, где $X_j = (x_{1j}, \dots, x_{nj})^T$, $j = 1, \dots, p$. На этой выборке известны результаты функции f , то есть имеем выборку $Y = (y_1, \dots, y_n)^T$ из случайной величины η и выборку $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ из случайной величины ϵ .

Пусть $x_i = (x_{i1}, \dots, x_{ip})$, тогда запишем зависимость f следующим образом:

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

Хотим по имеющейся выборке X оценить функцию f .

1.3 Линейная регрессия и МНК

Будем называть выборку $(x_1, \dots, x_n)^T$, участвующую в оценке функции f , обучающей выборкой и при $i = 1, \dots, n$

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

Будем называть выборку $(x'_1, \dots, x'_n)^T$, не участвующую в оценке функции f , тестовой выборкой и при $i = 1, \dots, n$

$$y'_i = f(x'_i) + \varepsilon'_i.$$

Далее будем предполагать, что рассматриваемая зависимость — линейная. Считаем, что $Y = (y_1, \dots, y_n)^T$ и $X = [X_1, \dots, X_p]$ — центрированы. Тогда модель многомерной линейной регрессии можно записать следующим образом:

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

где $\beta = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$ — параметры функции f .

Процесс подбора оптимального параметра β по обучающей выборке X называют обучением. Для оценки параметра β будем минимизировать среднюю квадратичную ошибку (МНК)

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

Решением метода наименьших квадратов (МНК) является вектор

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

и оценкой функции f , когда f — линейная функция, является

$$\hat{f}(x_i) = \sum_{j=1}^p \hat{\beta}_j x_{ij}.$$

1.4 Проблема МНК

Метод наименьших квадратов минимизирует среднюю квадратичную ошибку для обучающей выборки, но это не гарантирует, что на тестовой выборке \hat{f} будет также минимизировать среднюю квадратичную ошибку

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

Когда оценка \hat{f} дает ошибку больше на новых объектах, не участвовавших в обучении, чем на обучающей выборке ($\text{MSE}_{\text{test}} \gg \text{MSE}_{\text{training}}$) говорят, что происходит переобучение.

По мере увеличения количества параметров модели значение $\text{MSE}_{\text{training}}$ будет уменьшаться, а MSE_{test} нет. Это происходит потому что метод обучения будет сильно усложняться и пытаться объяснить то, что вызвано случайностью, а не свойствами функции f .

Рассмотрим от чего зависит MSE_{test} на генеральном языке. Пусть

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

Тогда математическое ожидание квадрата ошибки для x'_i

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

Таким образом, ошибка для объекта, не участвовавшего в обучении, зависит от дисперсии оценки функции f и квадрата смещения. Дисперсия оценки определяет то количество, на которое изменится \hat{f} , если бы мы получали эту оценку с использованием другого набора данных. Мы хотим, чтобы оценка \hat{f} не менялась сильно на разных обучающих выборках. Смещение \hat{f} характеризует ошибку, возникающую при аппроксимации реальной сложной функции f более простой моделью.

Как правило, при увеличении сложности метода (увеличение числа параметров) дисперсия будет увеличиваться, а смещение будет уменьшаться. За счет введения в оценку небольшого смещения, можно значительно уменьшить ее дисперсию и тем самым уменьшить MSE_{test} .

В случае многомерной линейной регрессии оценка \hat{f} не имеет смещения, так как нет смещения у оценки $\hat{\beta}$. Но чем больше дисперсия оценки $\hat{\beta}$

$$\text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1},$$

тем больше дисперсия \hat{f} . Когда матрица X близка к вырожденной, дисперсия $\hat{\beta}$ становится большой и MSE_{test} увеличивается.

2 Гребневая регрессия

2.1 Задача гребневой регрессии

Чтобы ввести смещение оценке параметра β добавим в $\text{MSE}_{\text{training}}$ множитель, который будет ограничивать параметр β . Тогда задача минимизации будет выглядеть следующим образом:

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

где $\lambda \geq 0$ — неотрицательный параметр регуляризации (tuning parameter). Такая модель минимизирует среднюю квадратичную ошибку, но при этом второй множитель $\lambda \sum_{j=1}^p \beta_j^2$ мал, когда параметры β_1, \dots, β_p близки к нулю.

При $\lambda = 0$, то гребневая регрессия совпадает с обычной регрессией, но при $\lambda \rightarrow \infty$ коэффициенты регрессии стремятся к нулю. Таким образом, для каждого значения λ будем получать свой оптимальный набор параметров β_1, \dots, β_p , поэтому для уменьшения MSE_{test} важно подобрать хорошее значение параметра λ .

2.2 Решение задачи минимизации и выбор параметра регуляризации

Решение гребневой регрессии можно выписать в явном виде:

$$\hat{\beta}_{\lambda}^R = (X^T X + \lambda I_p)^{-1} X^T Y. \quad (2.1)$$

Запишем модифицированное МНК решение гребневой регрессии через сингулярное разложение матрицы $X = VDU^T$, где матрицы V размера $n \times p$ и U размера $p \times p$ — ортогональные, а матрица $D = \text{diag}(\lambda_1, \dots, \lambda_p)$ — диагональная, на диагонали которой стоят собственные числа матрицы $X^T X$.

Тогда решение МНК:

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

Решение гребневой регрессии:

$$\hat{\beta}_{\lambda}^R = U(D^2 + \lambda I_p)^{-1} D V^T y = \sum_{j=1}^p \frac{\sqrt{\lambda_j}}{\lambda_j + \lambda} U_j (V_j^T Y). \quad (2.2)$$

При этом оценка функции f для выборки X запишется следующим образом:

$$X \hat{\beta}_{\lambda}^R = VDU^T \hat{\beta}_{\lambda}^R = V \text{diag}\left(\frac{\lambda_j}{\lambda_j + \lambda}\right) V^T Y = \sum_{j=1}^p \frac{\lambda_j}{\lambda_j + \lambda} V_j (V_j^T Y). \quad (2.3)$$

Такой способ записи решения задачи минимизации удобен для выбора параметра λ . Выражение (2.1) можно переписать, используя (2.2) и (2.3). Таким образом, необходимо один раз произвести сингулярное разложение матрицы X , а затем несложным образом вычислять функционал (2.1) для интересующих значений параметра λ .

Как выбрать параметр λ :

1. выбираем сетку значений λ ;

2. вычисляем ошибку кросс-проверки для каждого значения λ , используя запись (2.1) через сингулярное разложение;
3. выбираем λ с наименьшим значением ошибки кросс-проверки;
4. перестраиваем модель со всеми наблюдениями с выбранным значением λ .

2.3 Свойства и проблемы

Свойства:

1. Рассмотрим вероятностную интерпретацию гребневой регрессии. Предполагая, что
 - $\beta = (\beta_1, \dots, \beta_p)^T$ имеет априорное распределение $p(\beta)$;
 - $f(Y|X, \beta)$ — функция правдоподобия исходных данных;
 - линейная модель имеет независимые и нормально распределенные ошибки;
 - $p(\beta) = \prod_{j=1}^p g(\beta_j)$ для некоторой плотности g .

По теореме Байеса при фиксированном X апостериорное распределение $p(\beta|X, Y)$ пропорционально $f(Y|X, \beta)p(\beta|X) = f(Y|X, \beta)p(\beta)$ и если g — плотность $N(0, \lambda)$, то оценка апостериорного максимума β совпадает с решением гребневой регрессии.

2. Так как решение гребневой регрессии можно выписать в явном виде (2.2), то покажем, что оценка β имеет смещение. Пусть $X^T X = \Sigma$ и $Y = (y_1, \dots, y_n)^T$, тогда

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

Посчитаем математическое ожидание

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

При $\lambda = 0$ оценка совпадает с МНК оценкой и не имеет смещения.

Проблемы:

1. Оценки МНК инварианты относительно умножения признака на константу, то есть значение $f(x_j) \hat{\beta}_j$ не зависит от масштаба j -го признака. В случае гребневой регрессии инвариант относительно масштаба теряется, оценки МНК гребневой регрессии могут сильно измениться при умножении заданного признака на константу. Поэтому гребневую регрессию нужно использовать после стандартизации признаков.

2. В конечную модель входят все начальные признаки, если признаков много, то усложняется интерпретация.

3 Лассо

3.1 Задача минимизации Лассо

Метод Лассо, как и гребневая регрессия, вводит смещение у оценки β , при этом в конечную модель Лассо могут входить не все начальные признаки, что решает проблему гребневой регрессии. Это происходит за счет другого множителя в задачи минимизации, который регулирует коэффициенты β_1, \dots, β_p , однако из-за этого усложняется решение задачи. Задача минимизации Лассо:

$$\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j| \rightarrow \min_{\beta}, \quad (3.1)$$

где $\lambda \geq 0$ — неотрицательный параметр регуляризации (tuning parameter). Как и в гребневой регрессии $\lambda \sum_{j=1}^p |\beta_j|$ мало, когда β_1, \dots, β_p близки к нулю. При увеличении параметра λ некоторые коэффициенты регрессии становятся равными нулю. Как и в гребневой регрессии необходимо выбрать хорошее значение λ .

3.2 Решение задачи минимизации и выбор параметра регуляризации

Задача минимизации (3.1) 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 λ соответствует параметр s . Чем меньше s , тем больше нулевых значений коэффициентов β . Такую задачу можно решить, так как это задача квадратичного программирования.

Значение параметра λ выбирается как в гребневой регрессии с помощью кросс-проверки.

3.3 Свойства

Свойства:

1. Рассмотрим вероятностную интерпретацию Лассо. Предполагая, что

- $\beta = (\beta_1, \dots, \beta_p)^T$ имеет априорное распределение $p(\beta)$;
- $f(Y|X, \beta)$ — функция правдоподобия исходных данных;
- линейная модель имеет независимые и нормально распределенные ошибки;
- $p(\beta) = \prod_{j=1}^p g(\beta_j)$ для некоторой плотности g .

По теореме Байеса при фиксированном X апостериорное распределение $p(\beta|X, Y)$ пропорционально $f(Y|X, \beta)p(\beta|X) = f(Y|X, \beta)p(\beta)$ и если g — плотность распределения Лапласа с нулевым средним и параметром масштаба λ , то оценка апостериорного максимума β является решением Лассо.

2. Лассо обнуляет некоторые коэффициенты β_1, \dots, β_p . Задачу минимизации гребневой регрессии можно также рассмотреть в виде

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

Рассмотрим простой случай, когда $p = 2$. Тогда выражение $\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2$ будет представлять собой эллипс, центр которого находится в точке $\hat{\beta}$ (см. рис. 1). Предположим, что центр эллипса не удовлетворяет ограничениям $\sum_{j=1}^p \beta_j^2 \leq s$ и $\sum_{j=1}^p |\beta_j| \leq s$, то есть лежит вне круга в случае гребневой регрессии и вне ромба в случае Лассо. Тогда решения задач минимизации будут лежать на границе возможных значений. Данная ситуация изображена на рис. 1, можно заметить, что для Лассо существует гораздо больше различных эллипсов, которые пересекались бы с ромбом (ограничениями) таким образом, чтобы один из коэффициентов был равен нулю.

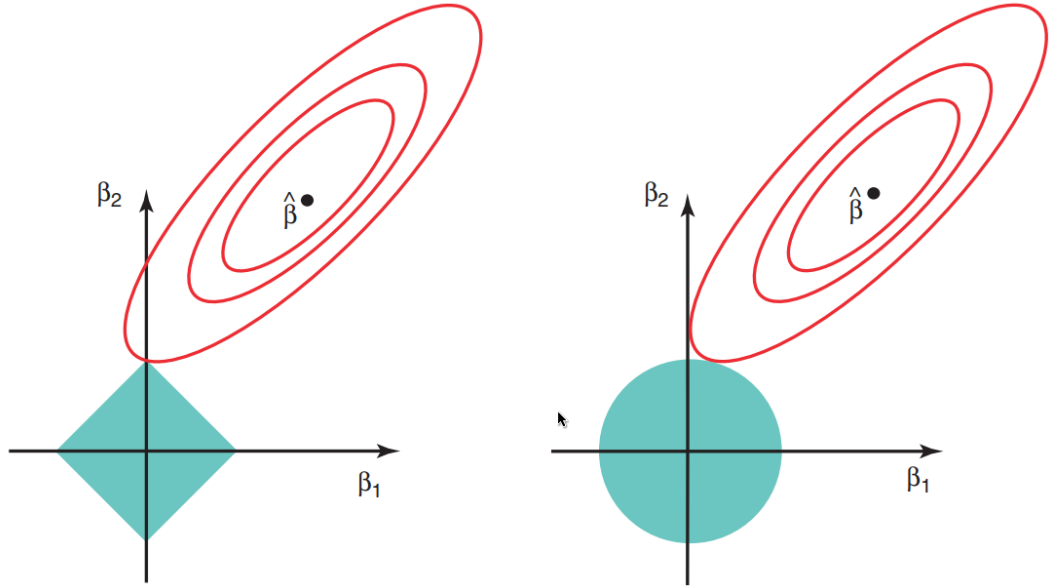


Рис. 1. Границы ошибки $\sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2$ и ограничений $\sum_{j=1}^p |\beta_j| \leq s$ для Лассо (слева) и $\sum_{j=1}^p \beta_j^2 \leq s$ для гребневой регрессии (справа).

4 Сравнение гребневой регрессии и Лассо

Рассмотрим простой случай, когда $n = p$ и X — диагональная матрица с 1 на диагонали. Тогда задача минимизации для обычной регрессии принимает вид

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

Решение МНК:

$$\hat{\beta}_j = y_j.$$

Решение гребневой регрессии будет принимать вид:

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

Решение Лассо:

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

Рисунок 4 иллюстрирует разный подход к регулированию коэффициентов β_1, \dots, β_p . Гребневая регрессия уменьшает каждый коэффициент с равной пропорцией. Лассо уменьшает значения коэффициентов на одинаковое значение, при этом если коэффициент по модулю меньше $\lambda/2$, то его значение становится равным нулю. В случае более общей матрицы X ситуация немного сложнее, но принцип регуляризации сохраняется.

В целом нельзя выделить ни одну из моделей (Лассо или гребневая регрессия) как лучшую. С помощью кросс-проверки можно определить какой подход лучше.

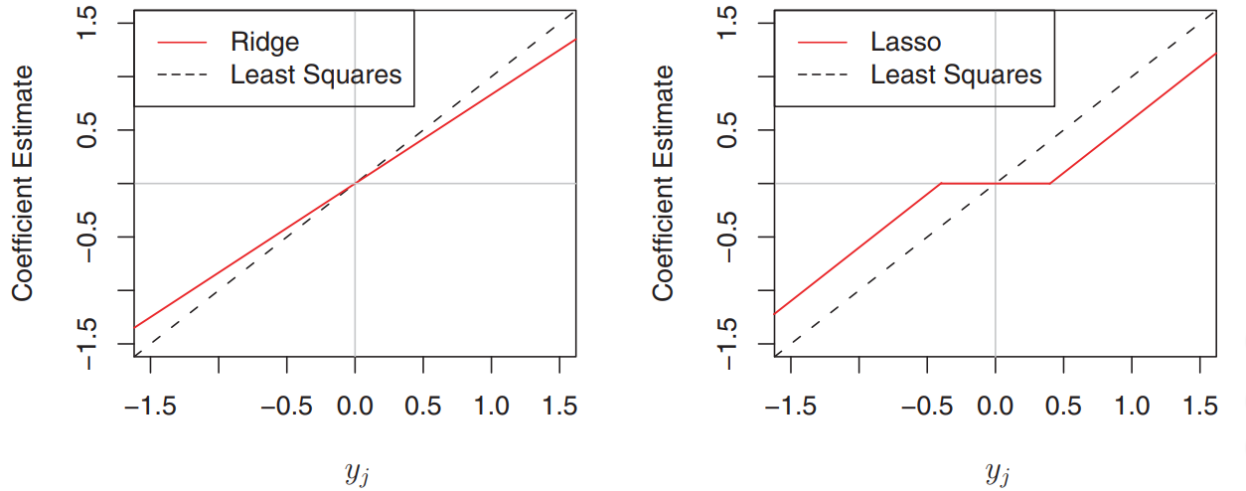


Рис. 2. Случай $n = p$. Слева сравнение МНК и гребневой регрессии. Справа сравнение МНК и Лассо.