

L_1 -регуляризация линейной регрессии. Регрессия наименьших углов (алгоритм LARS).

L_p -регуляризация линейной регрессии

Рассмотрим классическую модель линейной регрессии:

$$t = \sum_{j=1}^d w_j x(j) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

Поиск весов \mathbf{w} с помощью максимизации правдоподобия выборки в этой модели эквивалентен методу наименьших квадратов:

$$\sum_{n=1}^N \left[t_n - \sum_{j=1}^d w_j x_n(j) \right]^2 = \|\mathbf{t} - w_1 \mathbf{x}_1 - \dots - w_d \mathbf{x}_d\|^2 = \|\mathbf{t} - X\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}}. \quad (1)$$

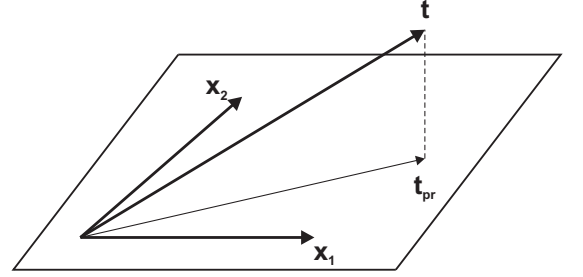


Рис. 1: Графическая иллюстрация линейной регрессии.

Здесь $\mathbf{x}_i \in \mathbb{R}^N$ — значения i -ого признака для всех объектов в выборке, $X = [\mathbf{x}_1, \dots, \mathbf{x}_d]$. Заметим, что введенное обозначение \mathbf{x}_i отличается от стандартного, когда под \mathbf{x}_i имеется ввиду i -ый объект выборки. Здесь и далее будем предполагать, что выборка является нормализованной, т.е. $\sum_{n=1}^N t_n = 0$, $\sum_{n=1}^N x_i(n) = 0$, $\frac{1}{N} \mathbf{t}^T \mathbf{t} = 1$, $\frac{1}{N} \mathbf{x}_i^T \mathbf{x}_i = 1$.

Задача (1) имеет простую геометрическую интерпретацию — поиск проекции вектора \mathbf{t} на гиперплоскость с направляющими векторами $\mathbf{x}_1, \dots, \mathbf{x}_d$ (см. рис. 1). Эта задача может быть решена аналитически:

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{t}, \quad \mathbf{t}_{pr} = X\mathbf{w} = X(X^T X)^{-1} X^T \mathbf{t}. \quad (2)$$

Данное решение для \mathbf{w} соответствует псевдорешению системы линейных уравнений $X\mathbf{w} = \mathbf{t}$.

Для того, чтобы избежать переобучения линейной регрессии, необходимо наложить ограничения на вариабельность решающего правила. Это можно сделать путем ограничения нормы вектора весов \mathbf{w} :

$$\begin{aligned} \|\mathbf{t} - X\mathbf{w}\|^2 &\rightarrow \min_{\mathbf{w}}, \\ \|\mathbf{w}\|_{L_p}^p &\leq b. \end{aligned} \quad (3)$$

Здесь под символом $\|\mathbf{w}\|_{L_p}$ понимается L_p норма вида $\sqrt[p]{\sum_{j=1}^d |w_j|^p}$. Традиционно вместо задачи оптимизации (3) рассматривают задачу оптимизации регуляризованного функционала:

$$\|\mathbf{t} - X\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|_{L_p}^p \rightarrow \min_{\mathbf{w}}, \quad \lambda \geq 0. \quad (4)$$

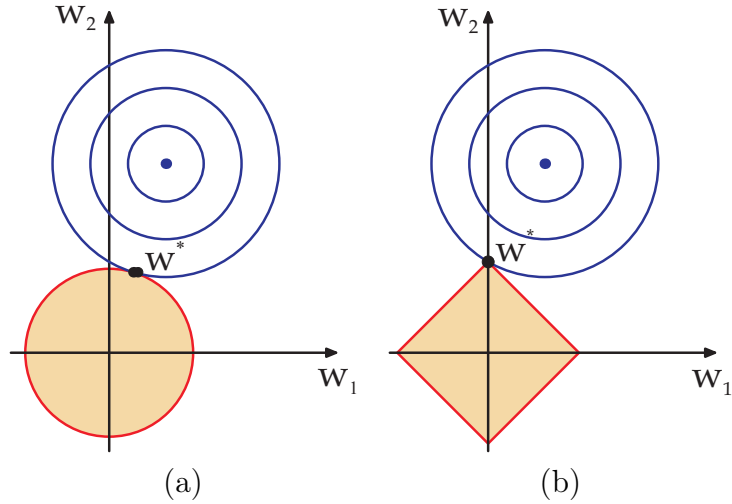


Рис. 2: Контуры минимизируемого функционала (синие кривые) и ограничения на норму весов. Случай (а) соответствует квадратичной норме, случай (b) — L_1 норме.

Нетрудно показать, что задачи оптимизации (3) и (4) эквивалентны при $p \geq 1$, т.е. когда все рассматриваемые функции являются выпуклыми. Введем функцию Лагранжа $L(\mathbf{w}, \lambda) = \|\mathbf{t} - X\mathbf{w}\|^2 + \lambda(\|\mathbf{w}\|_{L_p}^p - b)$. Тогда согласно выпуклому варианту теоремы Куна-Таккера необходимым и достаточным условием существования решения $\hat{\mathbf{w}}$ в задаче (3) является наличие $\lambda \geq 0$:

1. *Принцип минимума*: $L(\hat{\mathbf{w}}, \lambda) = \min_{\mathbf{w}} L(\mathbf{w}, \lambda)$,
2. *Условие дополняющей нежесткости*: $\lambda(\|\hat{\mathbf{w}}\|_{L_p}^p - b) = 0$.

Заметим, что для достаточности условий 1 и 2 требуется также выполнение условия Слейтера, т.е. существования $\mathbf{w} : \|\mathbf{w}\|_{L_p}^p < b$. Очевидно, что это условие выполнено при $b > 0$. Задача оптимизации (4) эквивалентна условию 1. Рассмотрим выполнение условия 2. Это условие эквивалентно наступлению одного из двух событий: $\lambda = 0$ или $\|\mathbf{w}\|_{L_p}^p = b$. Если $\lambda = 0$, то оптимальная точка $\hat{\mathbf{w}}$ лежит внутри области $\|\mathbf{w}\|_{L_p}^p < b$. Следовательно, ограничение $\|\mathbf{w}\|_{L_p}^p \leq b$ становится излишним, и задача оптимизации (3) переходит в задачу оптимизации без ограничений, что эквивалентно задаче (4) при $\lambda = 0$. Пусть $\lambda > 0$, $\|\mathbf{w}\|_{L_p}^p = b$. Выполнения этого условия легко добиться в задаче (4), просто обозначив через b значение нормы весов, оптимальных с точки зрения задачи (4).

Рассмотрим оптимальные решения задачи (3) при различных p . Можно показать, что в случае $p \leq 1$ оптимальное решение $\hat{\mathbf{w}}$ обладает свойством разреженности, т.е. часть весов в точности равна нулю. В случае $p > 1$ строго нулевые веса в оптимальном решении практически невозможны. Это свойство иллюстрируется на рис. 2. Заметим, что ситуация $p = 1$ является выделенной, т.к. в этом случае оптимизируемый функционал (4) является выпуклым, а оптимальное решение обладает свойством разреженности. Метод настройки весов линейной регрессии с помощью решения задачи (3) или (4) с L_1 нормой получил название LASSO (сокращение от Least Absolute Shrinkage and Selection Operator) [1]. Именно этот метод мы и будем рассматривать ниже.

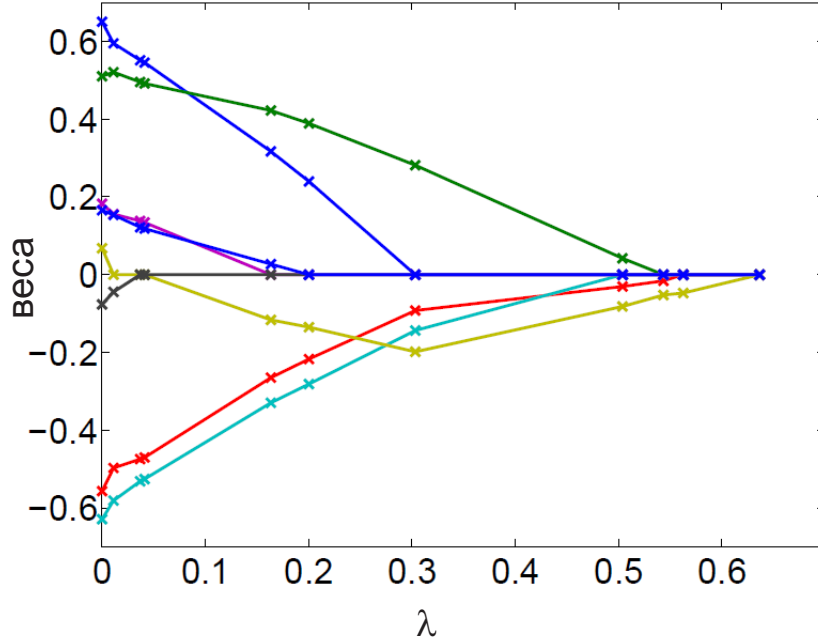


Рис. 3: «Путь регуляризации» для LASSO.

Теоретические свойства LASSO

Рассмотрим т.н. «путь регуляризации», т.е. совокупность решений $\hat{\mathbf{w}}$ задачи (4) для всех значений $\lambda \geq 0$. Можно показать, что этот путь является кусочно-линейным. Пусть при фиксированном λ нам известны знаки весов \mathbf{s} в оптимальном решении $\hat{\mathbf{w}}$:

$$\begin{aligned} s_j &= 0, \text{ если } \hat{w}_j = 0, \\ s_j &= 1, \text{ если } \hat{w}_j > 0, \\ s_j &= -1, \text{ если } \hat{w}_j < 0. \end{aligned}$$

В этом случае задача оптимизации (4) переходит в следующую:

$$\|\mathbf{t} - X\mathbf{w}\|^2 + \lambda \mathbf{s}^T \mathbf{w} \rightarrow \min_{\mathbf{w}}.$$

Данная задача может быть решена аналитически:

$$\begin{aligned} \hat{\mathbf{w}}_J(\lambda, \mathbf{s}) &= (X_J^T X_J)^{-1} (X_J^T \mathbf{t} - \frac{1}{2} \lambda \mathbf{s}), \\ \hat{\mathbf{w}}_{J^c}(\lambda, \mathbf{s}) &= \mathbf{0}. \end{aligned} \tag{5}$$

Здесь $J = \{j : s_j \neq 0\}$, $J^c = \{j : s_j = 0\}$, $X_J = X_{:,J}$. Можно показать, что необходимыми и достаточными условиями оптимальности решения (5) являются следующие:

$$\begin{aligned} \text{sign}(\hat{\mathbf{w}}_J(\lambda, \mathbf{s})) &= \mathbf{s}_J, \\ \|X_{J^c}^T (X_J \hat{\mathbf{w}}_J(\lambda, \mathbf{s}) - \mathbf{t})\|_{\infty} &\leq \lambda. \end{aligned} \tag{6}$$

Решение (5) является линейной функцией относительно λ . Это решение в совокупности с условиями (6) показывают, что итоговый путь регуляризации является кусочно-линейным (см.

рис. 3). Таким образом, возможный алгоритм решения задачи (4) состоит в поиске точек изменения направления на этом пути.

Существует ряд теоретических результатов относительно способности метода LASSO находить истинные релевантные веса (свойство состоятельности). Пусть наблюдаемая выборка данных сгенерирована из вероятностной модели линейной регрессии с весами \mathbf{w}^{true} , где часть весов равна нулю. Назовем метод LASSO состоятельным по знакам, если при стремлении объема выборки N к бесконечности верно

$$\lim_{N \rightarrow +\infty} P(\exists \lambda \geq 0 : \text{sign}(\hat{\mathbf{w}}(\lambda)) = \text{sign}(\mathbf{w}^{true})) = 1.$$

Здесь через $\hat{\mathbf{w}}(\lambda)$ обозначены веса, полученные с помощью LASSO с коэффициентом регуляризации λ . Можно показать [3], что метод LASSO является состоятельным по знакам тогда и только тогда, когда при стремлении объема выборки к бесконечности выполняется следующее условие:

$$\|Q_{J^c J} Q_{J J}^{-1} \text{sign}(\mathbf{w}_J^{true})\|_{\infty} \leq 1. \quad (7)$$

Здесь $Q = \lim_{N \rightarrow +\infty} \frac{1}{N} X_N^T X_N$ – выборочная матрица ковариации при стремлении объема выборки к бесконечности, а $J = \{j : w_j^{true} \neq 0\}$. На практике множество J , а также знаки \mathbf{w}^{true} являются неизвестными. Поэтому условие (7) приходится проверять для всех возможных знаков весов и всех возможных множеств J .

Рассмотрим условие, при котором требование (7) гарантированно выполнено. Обозначим через q мощность множества J , т.е. количество истинных информативных весов. Пусть матрица ковариации Q_N для выборки объема N имеет единицы на диагонали и внедиагональные элементы q_{ij} . Пусть найдется константа $0 < c \leq 1$: $|q_{ij}| \leq \frac{c}{2q-1}$. Тогда условие (7) выполнено [3].

Результат (7) является теоретическим. Экспериментальные исследования показывают, что на практике LASSO не является состоятельным по знакам [4] и, как правило, отбирает больше информативных признаков, чем нужно [5]. Для более устойчивого поиска множества информативных признаков обычно генерируется набор бутстрапированных выборок, для каждой из которых находится набор информативных признаков с помощью LASSO. Затем итоговое множество информативных признаков образуется в результате пересечения всех множеств информативных признаков [4] или путем усреднения рангов полученных признаков [6] (так как LASSO строит кусочно-линейную траекторию в пространстве весов, то помимо набора информативных признаков данный метод выдает еще и порядок вхождения признаков в информативное множество).

Алгоритм LARS (регрессия наименьших углов)

Изначательно автор LASSO предлагал решать задачу оптимизации (4) путем сведения к выпуклой задаче квадратичного программирования [1]. Для этого представим каждый вес w_j в следующем виде: $w_j = w_j^+ - w_j^-$, $w_j^+, w_j^- \geq 0$. Это представление не единственно. Рассмотрим следующую задачу оптимизации:

$$\|\mathbf{t} - X(\mathbf{w}^+ - \mathbf{w}^-)\|^2 + \lambda \sum_{j=1}^d (w_j^+ + w_j^-) \rightarrow \min_{\mathbf{w}^+, \mathbf{w}^-},$$

$$w_j^+, w_j^- \geq 0, \quad j = 1, \dots, d.$$

Легко показать, что решение данной задачи квадратичного программирования совпадает с решением LASSO. Однако, применение общих методов решения задачи квадратичного

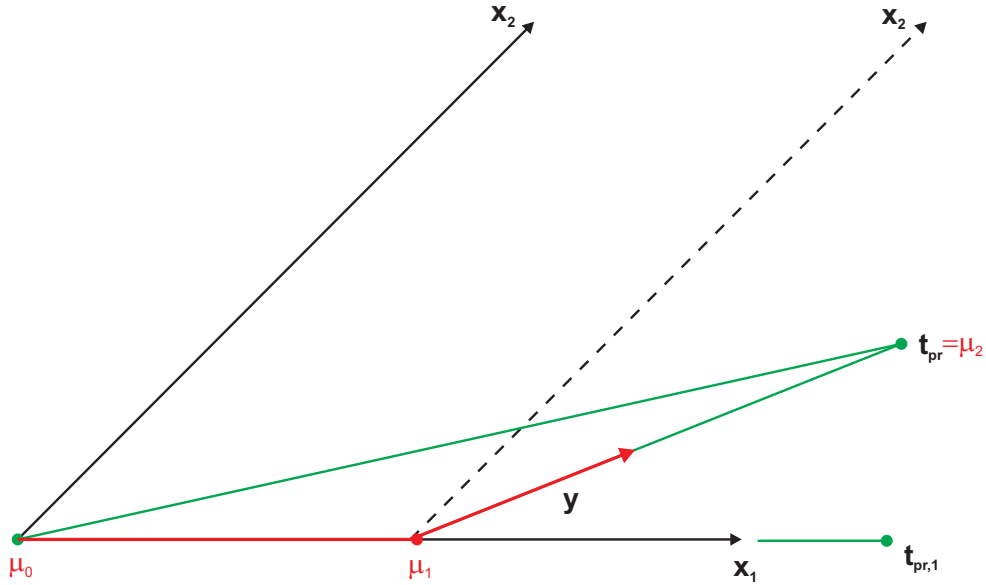


Рис. 4: Иллюстрация алгоритма LARS для двухмерного случая. В начале $\mu_0 = \mathbf{0}$ и наибольшая корреляция между $\mathbf{t} - \mu_0$ и векторами признаков достигается для \mathbf{x}_1 . На первом шаге алгоритма $\mu_1 = \mu_0 + \gamma_1 \mathbf{x}_1$, где γ_1 выбрано таким образом, чтобы корреляция между $\mathbf{t} - \mu_1$ и векторами признаков $\mathbf{x}_1, \mathbf{x}_2$ была бы одинакова. Далее выбирается вектор движения \mathbf{y} такой, вдоль которого корреляция между $\mathbf{t} - \mu$ и векторами признаков $\mathbf{x}_1, \mathbf{x}_2$ была бы все время одинаковой. Этот вектор соответствует биссектрисе угла, образуемый векторами $\mathbf{x}_1, \mathbf{x}_2$ с центром в μ_1 . На втором шаге $\mu_2 = \mu_1 + \gamma_2 \mathbf{y}$. Так как все имеющиеся признаки уже включены в текущий прогноз, то величина шага γ_2 выбирается в соответствии с величиной проекции $\mathbf{t} - \mu_1$ на вектор движения \mathbf{y} . Таким образом, μ_2 совпадает с оценкой наименьших квадратов \mathbf{t}_{pr} .

программирования требует значительного времени. Для метода LASSO был предложен специальный метод оптимизации [2], скорость работы которого сопоставима с поиском решения (2) (т.е. очень быстрый метод). Этот метод получил название LARS (сокращение от Least Angle Regression for laSso/Stepwise regression).

Алгоритм LARS направлен на построение оптимальной кусочно-линейной траектории в пространстве весов $\mathbf{w} \in \mathbb{R}^d$ (рис. 3) или прогнозов $\mu = X\mathbf{w} \in \mathbb{R}^N$. Таким образом, алгоритм LARS находит решения LASSO сразу для всех возможных значений коэффициента регуляризации λ .

Идея алгоритма LARS состоит в следующем (см. рис. 4). Будем обозначать через \mathbf{w}_j значение весов после j -ой итерации, а через $\mu_j = X\mathbf{w}_j$ значение прогноза для вектора \mathbf{t} после j -ой итерации. В начале $\mu_0 = \mathbf{0}$, $\mathbf{w}_0 = \mathbf{0}$. Найдём среди всех признаков \mathbf{x}_i наиболее коррелированный (по модулю) с \mathbf{t} . Обозначим его через \mathbf{x}_m , а знак корреляции через s_m . На первом шаге алгоритма модифицируем текущее значение прогноза μ_0 путем движения вдоль наиболее коррелированного вектора с учетом знака корреляции $\mu_1 = \mu_0 + \gamma s_m \mathbf{x}_m$, $\gamma > 0$. В результате такого движения величина корреляции между текущим остатком $\mathbf{t} - \mu_1$ и признаком \mathbf{x}_m будет уменьшаться по мере увеличения шага γ . Обозначим эту корреляцию через r_γ . Будем двигаться в выбранном направлении до тех пор, пока корреляция r_γ не совпадет с корреляцией между $\mathbf{t} - \mu_1$ и каким-нибудь еще признаком \mathbf{x}_k . Затем, на втором шаге алгоритма, будем модифицировать текущий прогноз μ_1 путем движения вдоль такого направления \mathbf{y}_1 , что корреляция между текущим остатком $\mathbf{t} - \mu_2$ и каждым из множества активных признаков $\mathbf{x}_m, \mathbf{x}_k$

была бы одинакова. Будем двигаться вдоль выбранного направления, пока еще один признак не окажется равно коррелированным с текущим остатком, что и признаки из активного множества. Тогда добавим этот признак в активное множество и найдем новое эквикоррелированное направление \mathbf{y}_2 , вдоль которого все признаки из активного множества являются одинаково коррелированными с текущим остатком. Алгоритм заканчивает работу, когда все имеющиеся (линейно независимые) признаки входят в активное множество. На последнем шаге алгоритма величина шага определяется значением корреляции между текущим остатком $\mathbf{t} - \boldsymbol{\mu}_{d-1}$ и последним эквикоррелированным направлением \mathbf{y}_{d-1} . При этом $\boldsymbol{\mu}_d = \mathbf{t}_{pr}$, где \mathbf{t}_{pr} является оценкой наименьших квадратов. Количество шагов алгоритма LARS определяется количеством признаков d .

Рассмотрим теперь шаги алгоритма LARS более подробно. Обозначим через $R = \frac{1}{N}X^T X$ выборочную матрицу ковариации выборки. Так как все вектора признаков \mathbf{x}_i считаются нормализованными, то значения в этой матрице R_{ij} определяют корреляцию между векторами \mathbf{x}_i и \mathbf{x}_j , а $R_{ii} = 1$. Обозначим через $\mathbf{r} = \frac{1}{N}X^T \mathbf{t}$ корреляцию между векторами признаков и целевым вектором \mathbf{t} .

Поиск оптимальной величины шага γ на первой итерации. Пусть на первом шаге алгоритма LARS $m = \arg \max_j |r_j|$, $s_m = \text{sign}(r_m)$. Найдем значение шага γ вдоль направления $s_m \mathbf{x}_m$. Корреляция между текущим остатком $\mathbf{t} - \gamma s_m \mathbf{x}_m$ и вектором $s_m \mathbf{x}_m$ составляет

$$\text{corr}(\mathbf{t} - \gamma s_m \mathbf{x}_m, s_m \mathbf{x}_m) = \frac{\frac{1}{N}(\mathbf{t} - \gamma s_m \mathbf{x}_m)^T s_m \mathbf{x}_m}{\sqrt{\mathbb{D}(\mathbf{t} - \gamma s_m \mathbf{x}_m)}} = \frac{s_m r_m - \gamma}{\sqrt{\mathbb{D}(\mathbf{t} - \gamma s_m \mathbf{x}_m)}}.$$

Здесь через $\mathbb{D}(\mathbf{t} - \gamma s_m \mathbf{x}_m)$ обозначена выборочная дисперсия случайной величины $t - \gamma s_m x_m$. Найдем корреляцию между текущим остатком и вектором \mathbf{x}_j , $j \neq m$:

$$\text{corr}(\mathbf{t} - \gamma s_m \mathbf{x}_m, \mathbf{x}_j) = \frac{r_j - \gamma s_m R_{mj}}{\sqrt{\mathbb{D}(\mathbf{t} - \gamma s_m \mathbf{x}_m)}}.$$

Приравнивая корреляции между собой, получаем:

$$\gamma_j^+ = \frac{s_m r_m - r_j}{1 - s_m R_{mj}}.$$

Вычисляя корреляцию между текущим остатком и вектором $-\mathbf{x}_j$, $j \neq m$, получаем:

$$\gamma_j^- = \frac{s_m r_m + r_j}{1 + s_m R_{mj}}.$$

Таким образом, оптимальная величина шага $\gamma = \min_{j \neq m} \gamma_j$, где $\gamma_j = \min^+(\gamma_j^+, \gamma_j^-)$. Под символом \min^+ понимается минимум среди положительных аргументов. После выполнения первого шага в активное множество признаков $\{s_m \mathbf{x}_m\}$ добавляется признак $s_k \mathbf{x}_k$, где $k = \arg \min_{j \neq m} \gamma_j$, $s_k = 1$, если $\gamma_k = \gamma_k^+$, $s_k = -1$ в противном случае.

Поиск эквикоррелированного направления \mathbf{y} . Пусть J – текущее активное множество признаков, $X_J^s = [\dots s_j \mathbf{x}_j \dots]$ – признаки из активного множества с учетом знаков \mathbf{s} , $\boldsymbol{\mu}$ – текущий прогноз. Тогда текущий эквикоррелированный вектор \mathbf{y} можно найти из следующих трех соображений. Во-первых, вектор \mathbf{y} является линейной комбинацией признаков из активного множества:

$$\mathbf{y} = X_J^s \mathbf{v}_J. \quad (8)$$

Здесь через \mathbf{v}_J обозначены веса этой линейной комбинации. Во-вторых, можно считать, что \mathbf{y} имеет единичную дисперсию, т.к. нас интересует только эквикоррелированное направление:

$$\frac{1}{N} \mathbf{y}^T \mathbf{y} = 1. \quad (9)$$

В-третьих, по определению эквикоррелированного направления \mathbf{y} должен иметь одинаковую корреляцию (обозначим ее через a) со всеми векторами признаков из активного множества, т.е.

$$\frac{1}{N} (X_j^s)^T \mathbf{y} = a \mathbf{1}_J. \quad (10)$$

Здесь через $\mathbf{1}_J$ обозначен вектор, состоящий из единиц. Подставляя (8) в (10), получаем

$$\frac{1}{N} (X_j^s)^T X_j^s \mathbf{v}_J = a \mathbf{1}_J. \quad (11)$$

Матрица $R_j^s = \frac{1}{N} (X_j^s)^T X_j^s$ является матрицей корреляций между векторами признаков из активного множества с учетом их знаков. Эта матрица может быть получена из подматрицы $R_{J,J}$ следующим образом: $R_j^s = D_J R_J D_J$, где $D_J = \text{diag}(\mathbf{s}_J)$. Теперь мы можем найти \mathbf{v}_J из условия (11):

$$\mathbf{v}_J = a (D_J R_J D_J)^{-1} \mathbf{1}_J. \quad (12)$$

Подставляя (8) в (9), получаем

$$\mathbf{v}_J^T \left[\frac{1}{N} (X_j^s)^T X_j^s \right] \mathbf{v}_J = 1.$$

Объединяя этот результат с (12), находим a :

$$a = \frac{1}{\sqrt{\mathbf{1}_J^T (D_J R_J D_J)^{-1} \mathbf{1}_J}}.$$

Подставляя найденное значение a в (12), находим \mathbf{v}_J , по которому в свою очередь из условия (8) находим \mathbf{y} . Обозначим через a_j корреляцию между \mathbf{y} и \mathbf{x}_j , $j \notin J$. Она может быть вычислена как

$$a_j = \frac{1}{N} \mathbf{x}_j^T \mathbf{y} = \mathbf{x}_j^T \frac{1}{N} X_j^s \mathbf{v}_J = (D_J R_{Jj})^T \mathbf{v}_J.$$

Поиск оптимальной величины шага γ вдоль эквикоррелированного направления \mathbf{y} .

Пусть $\boldsymbol{\mu}$ – текущий прогноз, \mathbf{y} – текущее эквикоррелированное направление. Нам необходимо вычислить величину шага γ вдоль эквикоррелированного направления для получения нового прогноза $\boldsymbol{\mu}^{new} = \boldsymbol{\mu} + \gamma \mathbf{y}$. Обозначим через r корреляцию признаков из активного множества с текущим остатком $\mathbf{t} - \boldsymbol{\mu}$, а через r_j – корреляцию остальных признаков с текущим остатком. Тогда

$$\begin{aligned} \forall m \in J : \text{corr}(\mathbf{t} - \boldsymbol{\mu} - \gamma \mathbf{y}, \mathbf{x}_m) &= \frac{r - \gamma a}{\sqrt{D(\mathbf{t} - \boldsymbol{\mu} - \gamma \mathbf{y})}}, \\ \forall j \notin J : \text{corr}(\mathbf{t} - \boldsymbol{\mu} - \gamma \mathbf{y}, \mathbf{x}_j) &= \frac{r_j - \gamma a_j}{\sqrt{D(\mathbf{t} - \boldsymbol{\mu} - \gamma \mathbf{y})}}, \\ \forall j \notin J : \text{corr}(\mathbf{t} - \boldsymbol{\mu} - \gamma \mathbf{y}, -\mathbf{x}_j) &= \frac{-r_j + \gamma a_j}{\sqrt{D(\mathbf{t} - \boldsymbol{\mu} - \gamma \mathbf{y})}}. \end{aligned}$$

Алгоритм 1 Алгоритм LARS

Вход: Коэффициенты корреляций между нормированными признаками $R \in \mathbb{R}^{d \times d}$, коэффициенты корреляций между нормированными признаками и нормированным целевым вектором $\mathbf{r} \in \mathbb{R}^d$.

Выход: Набор информативных признаков J в порядке убывания информативности, их знаки \mathbf{s} и веса (точки излома в пути регуляризации) $\mathbf{w}_1, \dots, \mathbf{w}_K$.

- 1: Найти максимальное количество линейно независимых признаков $d_{max} = \text{rank}(R)$.
 - 2: Инициализировать активное множество $J = \emptyset$, знаки признаков $s_J = \emptyset$ и веса $\mathbf{w} = \mathbf{0}$.
 - 3: Найти $m = \arg \max_j |r_j|$, $s_m = \text{sign}(r_m)$. Найти корреляцию признака m с текущим остатком $r_0 = s_m r_m$.
 - 4: **пока** $|J| < d_{max}$
 - 5: Добавить признак m в активное множество $J = J \cup \{m\}$, $s_J = s_J \cup \{s_m\}$.
 - 6: Найти $a = [\mathbf{1}_J^T (D_J R_J D_J)^{-1} \mathbf{1}_J]^{-1/2}$, где $R_J = R_{J,J}$, $D_J = \text{diag}(s_J)$, $\mathbf{1}_J$ – вектор из единиц.
 - 7: Найти $\mathbf{v}_J = a (D_J R_J D_J)^{-1} \mathbf{1}_J$, а также корреляцию между неактивными признаками и текущим эквикоррелированным направлением $a_j = (D_J R_{jJ})^T \mathbf{v}_J$, $j \in J^c$. // В том случае, если множество J состоит только из одного признака, то $a = 1$, $a_j = s_m R_{jm}$.
 - 8: Для $j \in J^c$ вычислить $\gamma_j^+ = (r_0 - r_j)/(a - a_j)$, $\gamma_j^- = (r_0 + r_j)/(a + a_j)$, $\gamma_j = \min^+\{\gamma_j^+, \gamma_j^-\}$. Найти $m = \arg \min_{j \in J^c} \gamma_j$. Если $\gamma_m = \gamma_m^+$, то $s_m = 1$, иначе $s_m = -1$.
 - 9: Найти новые веса $\mathbf{w}^{new} = \mathbf{w}^{old} + \gamma_m D_J \mathbf{v}_J$.
 - 10: Найти корреляцию признаков из активного множества с текущим остатком $r_0 = r_0 - \gamma_m a$ и корреляцию неактивных признаков с текущим остатком $r_j = r_j - \gamma_m a_j$.
 - 11: На последнем шаге веса совпадают с оценкой наименьших квадратов $\mathbf{w} = R^{-1} \mathbf{r}$.
-

Приравнивая соответствующие корреляции между собой, получаем:

$$\begin{aligned} \gamma &= \min_{j \notin J} \gamma_j, & \gamma_j &= \min^+(\gamma_j^+, \gamma_j^-), \\ \gamma_j^+ &= \frac{r - r_j}{a - a_j}, \\ \gamma_j^- &= \frac{r + r_j}{a + a_j}. \end{aligned}$$

Помимо величины шага γ мы находим также признак \mathbf{x}_k , который присоединяется к активному множеству, и его знак s_k : $k = \arg \min_{j \notin J} \gamma_j$, $s_k = 1$, если $\gamma_k = \gamma_k^+$, $s_k = -1$ в противном случае.

Остановка процесса оптимизации в алгоритме LARS. Очевидно, что в процессе оптимизации алгоритма LARS имеет смысл рассматривать только линейно независимые признаки. Только в этом случае матрица R_J будет невырожденной, а вычисления для a и \mathbf{v}_J будут корректными. Поэтому в алгоритме LARS множество активных признаков J наращивается до мощности d_{max} , где $d_{max} = \text{rank}(R)$. При этом на последнем шаге алгоритма оптимальные веса определяются через оценку наименьших квадратов $\mathbf{w} = (X_J^T X_J)^{-1} X_J^T \mathbf{t}$. Общая схема алгоритма LARS представлена в алгоритме 1.

Модификация LARS/LASSO. Алгоритм LARS, описанный выше, не решает задачу оптимизации для LASSO. Можно показать, что для LASSO справедливо следующее утверждение. Пусть $\boldsymbol{\mu} = X \mathbf{w}$ является решением LASSO для некоторого параметра

регуляризации λ . Тогда

$$\text{sign}(w_j) = s_j = \text{sign}(\mathbf{t} - \boldsymbol{\mu})^T \mathbf{x}_j, \quad j \in J,$$

т.е. знаки весов оптимального решения совпадают со знаками корреляций между текущим остатком и активными признаками. Метод LARS, вообще говоря, не гарантирует выполнения данного условия. Однако, можно предложить небольшую модификацию алгоритма LARS, которая будет решать задачу оптимизации для LASSO.

Пусть на текущем шаге алгоритма LARS вычислено \mathbf{v}_J и новые веса $\mathbf{w}^{new}(\gamma)$ определяются как $\mathbf{w}^{old} + \gamma D_J \mathbf{v}_J$. Таким образом, в точке $\gamma_j = -w_j^{old}/(s_j v_j)$ новый вес w_j^{new} меняет знак относительно w_j^{old} . Самое первое подобное изменение происходит для веса w_j^{new} для величины шага

$$\tilde{\gamma} = \min_j^+ \gamma_j, \quad \tilde{j} = \arg \min_j^+ \gamma_j.$$

Если все $\gamma_j < 0$, то положим по определению $\tilde{\gamma} = +\infty$. Тогда сделаем следующую модификацию алгоритма LARS:

$$\text{Если } \tilde{\gamma} < \gamma_m, \text{ то } \mathbf{w}^{new} = \mathbf{w}^{old} + \tilde{\gamma} D_J \mathbf{v}_J, \quad J = J \setminus \{\tilde{j}\}.$$

Здесь γ_m – оптимальная величина шага по эквикоррелированному направлению из алгоритма LARS. Можно показать [2], что с учетом данной модификации алгоритм LARS решает задачу оптимизации для LASSO сразу для всех значений коэффициента регуляризации λ .

Заметим, что в описанной модификации LARS/LASSO в некоторых ситуациях происходит уменьшение множества активных признаков. Подобная флуктуация мощности множества активных признаков увеличивает время работы алгоритма LARS/LASSO. В том случае, если нас интересует только само множество информативных признаков, а не значения весов оптимальной линейной комбинации, то лучше использовать исходный алгоритм LARS, который гарантированно сходится за d_{max} шагов.

Список литературы

- [1] *R. Tibshirani*. Regression Shrinkage and Selection via LASSO // Journal of the Royal Statistical Society, V. 58, I. 1, 1996, pp. 267–288.
- [2] *B. Efron, T. Hastie, I. Johnstone and R. Tibshirani*. Least Angle Regression // Annals of Statistics, V. 32, No. 2, 2004, pp. 407–499.
- [3] *P. Zhao, B. Yu*. On Model Selection Consistency of Lasso // Journal of Machine Learning Research, V. 7, 2006, pp. 2541–2563.
- [4] *F. Bach*. Bolasso: model consistent lasso estimation through the bootstrap // ICML, 2008.
- [5] *J. Lv, Y. Fan*. A unified approach to model selection and sparse recovery using regularized least squares // Annals of Statistics, V. 37, No. 6A, 2009, pp. 3498–3528.
- [6] *S. van Aelst, J.A. Khan and R.H. Zamar*. Robust Linear Model Selection Based on Least Angle Regression // Journal of the American Statistical Association, V. 102, 2007, pp. 1289–1299.