

# Анализ метода отбора признаков QPFS

Александр Адуенко

5 апреля 2021 г.

Рассматриваемый метод предложен в [1] и дается следующей оптимизационной задачей.

$$\begin{aligned} \frac{1}{2}(1 - \alpha)\mathbf{a}^\top \mathbf{Q}\mathbf{a} - \alpha\mathbf{b}^\top \mathbf{a} &\rightarrow \min_{\mathbf{a}} \\ \text{s.t. } \mathbf{a} &\geq 0, \sum_i a_i = 1. \end{aligned} \quad (1)$$

Решение задачи (1)  $\mathbf{a}^*$  определяет, какие признаки используются при построении модели. Признак  $j$  активен  $\iff a_j > 0$ .

Эта задача эквивалентна

$$\begin{aligned} \frac{1}{2}\mathbf{a}^\top \mathbf{Q}\mathbf{a} - \frac{\alpha}{1 - \alpha}\mathbf{b}^\top \mathbf{a} &\rightarrow \min_{\mathbf{a}} \\ \text{s.t. } \mathbf{a} &\geq 0, \sum_i a_i = 1. \end{aligned} \quad (2)$$

Обозначим эту задачу  $S\left(\underbrace{\frac{\alpha}{1 - \alpha}}_{\beta^{-1}}, 1\right)$ , где 1 указывает на норму решения. Рассмотрим задачу

$S(\beta^{-1}, \gamma)$  и сделаем замену переменной  $\mathbf{a} = \gamma\tilde{\mathbf{a}}$ , получим

$$\begin{aligned} \gamma^2 \left( \frac{1}{2}\tilde{\mathbf{a}}^\top \mathbf{Q}\tilde{\mathbf{a}} - \frac{1}{\beta\gamma}\mathbf{b}^\top \tilde{\mathbf{a}} \right) &\rightarrow \min_{\tilde{\mathbf{a}}} \\ \text{s.t. } \tilde{\mathbf{a}} &\geq 0, \sum_i \tilde{a}_i = 1, \end{aligned}$$

откуда задача  $S(\beta^{-1}, \gamma)$  эквивалентна задаче  $S((\beta\gamma)^{-1}, 1)$  в терминах активных признаков (сами же решения отличаются в  $\gamma$  раз), а потому мы имеем дело именно с однопараметрическим семейством и задание нормы решения, равной одному, не ограничивает общности.

Рассмотрим теперь замену переменной  $\mathbf{a} = \beta^{-1}\tilde{\mathbf{a}}$  в (2). Получим, что (2) эквивалентна при  $\alpha \in (0, 1)$  задаче

$$\begin{aligned} \frac{1}{2}\tilde{\mathbf{a}}^\top \mathbf{Q}\tilde{\mathbf{a}} - \mathbf{b}^\top \tilde{\mathbf{a}} &\rightarrow \min_{\tilde{\mathbf{a}}} \\ \text{s.t. } \tilde{\mathbf{a}} &\geq 0, \sum_i \tilde{a}_i = \beta. \end{aligned} \quad (3)$$

Наряду с задачей (3) можно рассмотреть задачу

$$\begin{aligned} \frac{1}{2}\tilde{\mathbf{a}}^\top \mathbf{Q}\tilde{\mathbf{a}} - \mathbf{b}^\top \tilde{\mathbf{a}} &\rightarrow \min_{\tilde{\mathbf{a}}} \\ \text{s.t. } \tilde{\mathbf{a}} &\geq 0, \sum_i \tilde{a}_i \leq \beta \end{aligned} \quad (4)$$

и соответствующую задачу без ограничения на норму вектора  $\tilde{\mathbf{a}}$

$$\begin{aligned} \frac{1}{2}\tilde{\mathbf{a}}^\top \mathbf{Q}\tilde{\mathbf{a}} - \mathbf{b}^\top \tilde{\mathbf{a}} &\rightarrow \min_{\tilde{\mathbf{a}}} \\ \text{s.t. } \tilde{\mathbf{a}} &\geq 0. \end{aligned} \quad (5)$$

Задача (1) не позволяет выбросить все признаки, поскольку  $\|\mathbf{a}\| = 1 > 0$ . Задача (4) эквивалентна ограничению неравенства в исходной задаче и, например, при  $\alpha = 0$  будет иметь решением исключение всех признаков. Далее приведем анализ свойств решения (1) с ограничением равенства и неравенства для разных значений  $\alpha$ .

Таблица 1. Свойства решения в методе QPFS в зависимости от параметра  $\alpha$

$\alpha$	Равенство	Неравенство
$\alpha = 0$	Используются все признаки, $\mathbf{Q}\mathbf{a}^* = \eta\mathbf{e}$ , $\eta > 0$	Выброшены все признаки
$\alpha \rightarrow 0$	Используются все признаки, $\mathbf{Q}\mathbf{a}^* \rightarrow \eta\mathbf{e}$ , $\eta > 0$	Решение (5)
$\alpha \rightarrow 1$	Сходимся к отбору одного признака с максимальным $b_j$	То же, что и в «равенство»
$\alpha = 1$	Отбор одного признака с максимальным $b_j$	То же, что и в «равенство»

В случае  $\alpha = 0$  для неравенства ( $\sum_i a_i \leq 1$ ) решением является  $\mathbf{a}^* = \mathbf{0}$ . В случае равенства ( $\sum_i a_i = 1$ ), чтобы минимизировать потери от необходимости иметь ненулевой  $\mathbf{a}$ , в оптимальном  $\mathbf{a}^*$  оптимизируемая функция  $\frac{1}{2}\mathbf{a}^\top \mathbf{Q}\mathbf{a}$  должна иметь одинаковый градиент по всем направлениям (так как иначе можно уменьшить одну координату, увеличить другую, оставив норму  $\mathbf{a}$  неизменной, уменьшив значение функции). Тот же результат можно получить и другим способом - из рассмотрения Лагранжиана задачи.

## Анализ решения QPFS в зависимости от параметра $\alpha$

Задачи (3) и (4) эквивалентны для некоторых  $\eta$  задаче

$$\frac{1}{2}\tilde{\mathbf{a}}^\top \mathbf{Q}\tilde{\mathbf{a}} - \mathbf{b}^\top \tilde{\mathbf{a}} + \eta \sum_j \mathbf{a}_j \rightarrow \min_{\tilde{\mathbf{a}}} \text{ s.t. } \tilde{\mathbf{a}} \geq 0,$$

что эквивалентно

$$\frac{1}{2}\tilde{\mathbf{a}}^\top \mathbf{Q}\tilde{\mathbf{a}} - \tilde{\mathbf{b}}^\top \tilde{\mathbf{a}} \rightarrow \min_{\tilde{\mathbf{a}}} \text{ s.t. } \tilde{\mathbf{a}} \geq 0,$$

где  $\tilde{b}_j = b_j - \eta$  и  $\eta$  монотонно убывает по  $\|\tilde{\mathbf{a}}\|_1 = \beta$ . При этом для случая неравенства  $\eta \geq 0$ , в для равенства  $\eta < 0$ , если  $\beta > \|\mathbf{a}^*\|_1$ , где  $\mathbf{a}^*$  есть решение задачи без ограничения на норму (5).

Таким образом, добавление ограничения на норму, фактически штрафует релевантность и происходит исключение тех признаков, у которых  $b_j < \eta$ , поскольку у них поправленная релевантность  $\tilde{b}_j$  становится отрицательной.

## Связь с лассо для линейной регрессии

Стандартная задача линейной регрессии с  $l_1$  регуляризацией (лассо) имеет вид

$$\frac{1}{2}\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \tau\|\mathbf{w}\|_1 \rightarrow \min_{\mathbf{w}}. \quad (6)$$

Предположим, что  $\mathbf{x}_j^\top \mathbf{x}_j = 1$ ,  $\mathbf{y}^\top \mathbf{y} = 1$ , то есть признаки и целевая переменная нормированы. В качестве функций сходства (similarity) и релевантности (relevance) в QPFS рассмотрим

корреляцию Пирсона (не модуль как обычно в QPFS!). Имеем

$$\frac{1}{2}\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \tau\|\mathbf{w}\|_1 = \frac{1}{2}\underbrace{\mathbf{y}^\top\mathbf{y}}_{=1} + \frac{1}{2}\mathbf{w}^\top \underbrace{\mathbf{X}^\top\mathbf{X}}_{\tilde{\mathbf{Q}}} \mathbf{w} - \underbrace{(\mathbf{X}^\top\mathbf{y})^\top}_{\tilde{\mathbf{b}}} \mathbf{w},$$

где учтена нормированность всех признаков и целевой переменной, откуда корреляция и ко-вариация совпадают;  $\tilde{\mathbf{Q}}, \tilde{\mathbf{b}}$  корреляции со знаком между признаками и признаками и целевой переменной соответственно. Отсюда задачу (6) можно переписать в виде

$$\frac{1}{2}\mathbf{w}^\top \tilde{\mathbf{Q}}\mathbf{w} - \tilde{\mathbf{b}}^\top \mathbf{w} + \tau\|\mathbf{w}\|_1 \rightarrow \min_{\mathbf{w}},$$

что может быть переписано в эквивалентную задачу с ограничением равенства (так же для неравенства) для некоторого  $\eta$

$$\begin{aligned} \frac{1}{2}\mathbf{w}^\top \tilde{\mathbf{Q}}\mathbf{w} - \tilde{\mathbf{b}}^\top \mathbf{w} &\rightarrow \min_{\mathbf{w}} \\ \text{s.t. } \|\mathbf{w}\|_1 &= \eta. \end{aligned} \tag{7}$$

Если все корреляции Пирсона между признаками и между признаками и целевой переменной неотрицательны (то есть векторы  $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{y}$  лежат в одном многомерном квадранте), то  $\tilde{\mathbf{Q}} = \mathbf{Q}, \tilde{\mathbf{b}} = \mathbf{b}$ , то есть задача (7) тождественна QPFS, но без ограничения  $\mathbf{w} \geq 0$ . Таким образом, QPFS можно рассматривать как lasso ограничением на неотрицательность весов, если все корреляции Пирсона между признаками и между признаками и целевой переменной неотрицательны. Далее, если истинный вектор весов в линейной регрессии  $\mathbf{w} \geq 0$ , то условие не неотрицательность оценки весов тоже избыточно (так как  $\mathbf{w}^*$  в задаче (7) и так будет неотрицательным, начиная с некоторого размера выборки), и QPFS будет полностью тождественен lasso. Таким образом, получаем условия тождественности QPFS( $\alpha$ ) методу lasso( $\tau$ ), где между  $\alpha$  и  $\tau$  существует некоторая связь:

- Rel = Sim = |Pearson correlation|;
- Нормированность признаков и целевой переменной:  $\mathbf{y}^\top\mathbf{y} = \mathbf{x}_j^\top\mathbf{x}_j = 1, j = 1, \dots, n$ ;
- Неотрицательность попарных корреляций:  $\mathbf{y}^\top\mathbf{x}_j \geq 0, \mathbf{x}_j^\top\mathbf{x}_l \geq 0, j, l = 1, \dots, n$ ;
- Неотрицательность истинного вектора весов  $\mathbf{w}^* \geq 0$ .

Отметим, что в lasso реализуется ситуация, когда выбрасываются все признаки, когда  $\tau \rightarrow \infty$ , что соответствует QPFS с ограничением неравенства (а не равенства) при  $\alpha = 0$ . Кроме того, «закритический» режим из QPFS с равенством, когда  $\|\mathbf{a}\|_1 = \beta > \|\mathbf{a}^*\|$ , где  $\mathbf{a}^*$  есть решение (5), в lasso не реализуется, поскольку, штраф за норму приводит всегда к ее сокращению, а потому  $\tau \rightarrow 0$  в lasso соответствует  $\beta \rightarrow \|\mathbf{a}^*\|$  в QPFS с ограничением типа равенства.

## Проблемы QPFS

Благодаря условию на неотрицательность коэффициентов  $\mathbf{a} \geq 0$  в задаче QPFS (3), штраф на  $\|\mathbf{a}\|_1$  становится штрафом на сумму коэффициентов, что упрощает оптимизацию. Авторы оригинального метода [1] прямо указывают на скорость оптимизации как на основное преимущество метода QPFS при сопоставимом качестве прогноза (например, с lasso) на тестовых выборках в рассмотренных наборах данных.

Как показано в предыдущей главе, при выполнении некоторых условий (в частности, неотрицательности истинного вектора параметров модели  $\mathbf{w}^*$ ), QPFS будет в точности эквивалентен методу lasso, то есть метод обладает лучшими свойствами с точки зрения оптимизации, при этом давая то же решение, что и lasso. Однако, когда условия эквивалентности не выполнены, метод начинает проигрывать lasso, поскольку не учитывает, например, что релевантность пары признаков может быть значительно выше, чем релевантность каждого из них.

**Пример.** Рассмотрим выборку  $(\mathbf{X}, \mathbf{y})$  в признаковом пространстве размерности  $n = 2$ .

$$x_{1i} \sim \mathcal{N}(x_{1i}|0, 1), y_i \sim \mathcal{N}(y_i|0, 1), x_{2i} = x_{1i} + \varepsilon y_i, \varepsilon > 0.$$

Истинная корреляция Пирсона первого признака и целевой переменной  $y$  равна 0, а корреляция со вторым – равна  $\varepsilon/\sqrt{1+\varepsilon^2}$ , схожесть двух признаков  $1/\sqrt{1+\varepsilon^2}$ . При малом  $\varepsilon$  выборочная корреляция обоих признаков с целевой переменной будет мала, а схожесть двух признаков – велика. При этом надежное восстановление целевой переменной  $y$  возможно только при наличии обоих признаков в выборке, что соответствует ситуации с отсутствием отбора признаков (малое  $\alpha$ ).

Добавим теперь в выборку  $N$  шумовых признаков. Истинное сходство каждого из таких признаков с целевой переменной равно 0 (такое же как и для признака 1) и с учетом того, что QPFS не учитывает взаимодействия между признаками, признаки 1 и 2 не имеют значительного преимущества по отношению к шумовым в терминах релевантности целевой переменной (признак 1 в точности шумовой в изоляции, так как независим от  $y$ ). При этом признаки 1 и 2 получают штраф за похожесть друг на друга. По этой причине при работе QPFS либо происходит исключение одного или обоих признаков 1 и 2 при исключении некоторых или всех шумовых, или оба признака 1 и 2 остаются, но вместе с ними остаются почти все или все шумы (см. эксперимент). В то же время Lasso учитывает взаимосвязи между признаками и не требует неотрицательности коэффициентов (ссылка на сравнение на этом датасете QPFS и lasso). В рассматриваемом примере не выполнено одно из условий эквивалентности lasso QPFS:

- $\mathbf{w}^* = (-1/\varepsilon, 1/\varepsilon)^T$  содержит отрицательные значения.

Подобный пример (нужно добавить в статью вместе с соответствующим экспериментом с  $N = 10$  или  $N = 100$  шумами):

1.  $x_{2i} = -x_{1i} + \varepsilon y_i \rightarrow$  не выполнено условие неотрицательной корреляции между признаками, а остальные условия эквивалентности выполнены.

**Замечание:** Подумай, можно ли построить пример, чтобы все условия кроме одного были выполнены для оставшихся условий эквивалентности QPFS и lasso.

## Стабильность модели

В статье Катруцы (ссылка) стабильность модели (выраженная, например, в терминах  $\lambda_{\max}(\mathbf{X}^T\mathbf{X})/\lambda_{\min}(\mathbf{X}^T\mathbf{X})$ ) имеет самостоятельную ценность и наряду с качеством прогноза на тестовой выборке определяет решение о превосходстве одного метода отбора признаков над другим. В этой статье мы предлагаем рассматривать стабильность как априорное знание, которое указывает на то, что априори мы считаем, что выборки с меньшим числом обусловленности на множестве активных признаков появляются чаще в рассматриваемой задаче, чем выборки с большим числом обусловленности. При отсутствии такого знания стабильность

стоит рассматривать в контексте повышения качества прогноза: если низкая стабильность модели ведет к снижению качества прогноза, стоит добавить штраф за низкую стабильность модели, если качество не снижается, а повышается при уменьшении стабильности, то не стоит предпочитать менее качественную, но более стабильную модель.

Обозначим  $\mathbf{X}(\mathbf{w})$  сужение матрицы признаков на множество признаков  $j : w_j \neq 0$ . Примером соответствующей задачи оптимизации, где есть априорное знание о том, что низкое число обусловленности более предпочтительно является

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \tau \lambda_{\max}(\mathbf{X}(\mathbf{w})^T \mathbf{X}(\mathbf{w})) / \lambda_{\min}(\mathbf{X}(\mathbf{w})^T \mathbf{X}(\mathbf{w})) \rightarrow \min_{\mathbf{w}},$$

что соответствует показательному распределению на число обусловленности активных признаков выборки с гиперпараметром  $\tau$ . В таком виде задача является тяжелой для оптимизации и требует перебора наборов активных признаков, например, с помощью генетического алгоритма.

**Пример 1.** Пусть  $y_i \sim \mathcal{N}(y_i|0, 1)$ ,  $x_{ij} = y_i + \nu \varepsilon_{ij}$ ,  $\varepsilon_{ij} \sim \mathcal{N}(\varepsilon_{ij}|0, 1)$ .

Оптимальная модель использует все признаки и осредняет их для получения наилучшего прогноза  $y_i$ :  $\mathbf{w}^* = (1/n, \dots, 1/n)^T$ . При этом число обусловленности  $\eta = \lambda_{\max}(\mathbf{X}^T \mathbf{X}) / \lambda_{\min}(\mathbf{X}^T \mathbf{X})$  при малом  $\nu$  может быть большим.

**Пример 2.** Пусть в выборке есть дубликат признака  $y_i \sim \mathcal{N}(y_i|0, 1)$ ,  $x_{i1}x_{i2} = y_i + \nu \varepsilon_i$ ,  $\varepsilon_i \sim \mathcal{N}(\varepsilon_i|0, 1)$ . В этом случае число обусловленности равно  $\infty$  и имеется неоднозначность решения, минимизирующего  $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ . Однако все эти решения имеют  $w_1 + w_2 = \text{const}$ , а потому если на тестовой выборке признаки 1 и 2 останутся идентичными, качество прогноза будет одинаковым независимо от того, какое разбиение этой константы между  $w_1$ ,  $w_2$ , мы предпочтем. Обычно для того, чтобы сделать решение однозначным, добавляют слабую квадратичную регуляризацию на  $\mathbf{w}$ , что из всех решения предпочитается то, где  $w_1^2 + w_2^2$  минимально, то есть  $w_1 = w_2 = \text{const}/2$ .

Заметим, что если есть основания полагать, что сильная мультиколлинеарность в обучающей выборке не будет продолжена в тестовой [2], то специальная обработка мультиколлинеарности приобретает важность, а конкретный вид поправок зависит от предположений об эволюции корреляций.

## Список литературы

- [1] *Rodriguez-Lujan I. et al.* Quadratic programming feature selection // Journal of Machine Learning Research. – 2010.
- [2] *Belsley D. A.* Collinearity and forecasting // Journal of Forecasting. – 1984. – Vol. 3. – No. 2. – P. 183-196.