

Об унимодальности непрерывного расширения критерия Акаике*

Ветров Д. П., Кропотов Д. А., Пташко Н. О.

vetrovd@yandex.ru, dmitry.kropotov@gmail.com, ptashko@inbox.ru

Москва, ВМиК МГУ, Вычислительный Центр РАН

В работе рассматривается применение непрерывного расширения критерия Акаике (САИС) к подбору параметров регуляризации в задаче обобщенной линейной регрессии. Значениями параметра регуляризации являются все симметричные неотрицательно определенные матрицы. Показывается, что на множестве всех таких матриц критерий САИС является унимодальным. Получено явное условие вырожденности решающего правила (нулевого решения). Показано, что данное условие остается справедливым для семейства диагональных неотрицательно определенных матриц, а также для семейства матриц, пропорциональных единичной.

Для решения задач выбора моделей широко применяется информационный критерий Акаике [1], предлагающий подход, основанный на теории информации. Несмотря на то, что этот метод изначально был предложен для выбора из конечного числа моделей, он может быть расширен на случай бесконечного семейства моделей. Такое расширение позволяет заменить полный перебор на конечном множестве направленным поиском максимума непрерывного функционала, что существенно ускоряет процедуру выбора. В работах [2, 3] предложено подобное обобщение информационного критерия Акаике (САИС) в применении для настройки параметров модели стационарной и нестационарной линейной регрессии. Одновременная настройка нескольких параметров модели с помощью критерия САИС обычно проводится градиентным или покоординатным методами. В этой связи возникает вопрос о наличии у непрерывного критерия Акаике локальных максимумов. В данной работе рассматривается задача выбора модели с помощью САИС в широком семействе всех симметричных неотрицательно определенных матриц. Доказано, что непрерывный критерий Акаике в этом семействе является унимодальной функцией. Выписано аналитическое решение для обобщенной линейной модели регрессии, в которой, согласно критерию Акаике, достигается наилучшая обобщающая способность. Также получено практически важное условие релевантности, допускающее непосредственную прямую проверку, которое позволяет отбросить заведомо неприемлемые модели. Рассматриваемое семейство матриц включает в себя важные подсемейства: семейство всех диагональных матриц и семейство матриц, пропорциональных единичной. Полученные здесь результаты для широкого семейства могут быть частично перенесены для данных подсемейств. В частности, показано, что условие релевантности остается справедливым и для рассматриваемых подсемейств.

*Работа выполнена при финансовой поддержке РФФИ, проекты № 08-01-00405, № 08-01-90016, № 08-01-90427, № 07-01-00211.

Непрерывное расширение критерия Акаике

Рассмотрим классическую задачу обобщенной линейной регрессии. Пусть $(X, \mathbf{t}) = \{(\mathbf{x}_i, t_i)\}_{i=1}^n$ — обучающая выборка, где $\mathbf{x}_i = (x_i^1, \dots, x_i^d) \in \mathbb{R}^d$ — вектор наблюдаемых признаков объекта, $t_i \in \mathbb{R}$ — значение зависимой переменной. Зафиксируем некоторое множество базисных функций $\{\varphi_i(\mathbf{x})\}_{i=1}^m$, $\varphi_j : \mathbb{R}^d \rightarrow \mathbb{R}$. Требуется найти вектор весов $\mathbf{w} \in \mathbb{R}$ такой, что функция

$$y(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\varphi}(\mathbf{x}) = \sum_{j=1}^m w_j \varphi_j(\mathbf{x})$$

приближала бы значения переменной t на объектах обучающей выборки X . Пусть $\Phi = (\varphi_{ij})_{n \times m} = (\varphi_j(\mathbf{x}_i))_{n \times m}$ — матрица базисных функций, вычисленных для каждого объекта обучающей выборки. Классический подход к обучению линейной регрессии состоит в оптимизации регуляризованного правдоподобия

$$\mathbf{w}_{\text{MP}} = \arg \max_{\mathbf{w}} p(\mathbf{t} | X, \mathbf{w}) p(\mathbf{w} | \alpha), \quad (1)$$

где

$$p(\mathbf{t} | X, \mathbf{w}) = \frac{1}{\sqrt{(2\pi)^n \sigma^n}} \exp\left(-\frac{1}{2\sigma^2} \|\Phi \mathbf{w} - \mathbf{t}\|^2\right) \quad (2)$$

— функция правдоподобия.

В качестве регуляризатора $p(\mathbf{w} | \alpha)$ часто рассматривается квадратичный функционал с некоторым коэффициентом регуляризации $\alpha \geq 0$:

$$p(\mathbf{w} | \alpha) = \left(\frac{\alpha}{2\pi}\right)^{m/2} \exp\left(-\frac{\alpha}{2} \sum_{j=1}^m w_j^2\right). \quad (3)$$

В методе релевантных векторов RVM [4] для автоматического отбора релевантных базисных функций семейство регуляризаторов (3) предлагается расширить, и для каждого веса w_j ввести свой

коэффициент регуляризации α_j :

$$\begin{aligned} p(\mathbf{w} | \boldsymbol{\alpha}) &= \prod_{j=1}^m \sqrt{\frac{\alpha_j}{2\pi}} \exp\left(-\frac{\alpha_j}{2} w_j^2\right) = \\ &= \sqrt{\frac{\det(R)}{(2\pi)^m}} \exp\left(-\frac{1}{2} \mathbf{w}^T R \mathbf{w}\right), \end{aligned} \quad (4)$$

где $R = \text{diag}(\alpha_1, \dots, \alpha_m)$ — матрица регуляризации, $\alpha_j \geq 0$. Расширим используемое в RVM семейство регуляризаторов на случай всех (необязательно диагональных) симметричных неотрицательно определенных матриц $R = R^T \succeq 0$. Такое семейство матриц позволяет не только находить релевантное подмножество базисных функций, но и одновременно выделять релевантные линейные комбинации исходных базисных функций.

Подбор матрицы регуляризации будем осуществлять, используя непрерывное расширение критерия Акаике, САИС [2]:

$$\begin{aligned} R &= \arg \max f(R); \\ f(R) &= \log p(\mathbf{t} | X, \mathbf{w}_{\text{MP}}) - \text{tr}(H(H+R)^{-1}). \end{aligned} \quad (5)$$

Здесь $H = -\nabla \nabla \log p(\mathbf{t} | X, \mathbf{w}) = \sigma^{-2} \Phi^T \Phi$.

Параметр σ также может быть найден путем максимизации САИС, что приводит к следующему итеративному процессу его вычисления на основе текущего значения R :

$$(\sigma^2)^{\text{new}} = \frac{\|\mathbf{t} - \Phi \mathbf{w}_{\text{MP}}\|^2}{n - \text{tr} H(H+R)^{-1} R(H+R)^{-1}}. \quad (6)$$

Решение задачи оптимизации

Обозначим через \mathbf{w}_{ML} оценку максимального правдоподобия на выборке X . Можно показать, что при условии (2) и (4) максимизация САИС (5) эквивалентна следующей оптимизационной задаче:

$$\begin{cases} -\frac{1}{2} \mathbf{w}_{\text{ML}}^T H(H+R)^{-1} H(H+R)^{-1} H \mathbf{w}_{\text{ML}} + \\ + \mathbf{w}_{\text{ML}}^T H(H+R)^{-1} H \mathbf{w}_{\text{ML}} - \\ - \text{tr}(H(H+R)^{-1}) \rightarrow \max_R; \\ R = R^T \succeq 0. \end{cases} \quad (7)$$

Обозначим $\mathbf{v} = H^{\frac{1}{2}} \mathbf{w}_{\text{ML}}$,

$$A = H^{\frac{1}{2}} (H+R)^{-1} H^{\frac{1}{2}} = (I + H^{-\frac{1}{2}} R H^{-\frac{1}{2}})^{-1}.$$

Тогда задача (7) может быть переписана в следующем виде:

$$\begin{cases} -\frac{1}{2} \mathbf{v}^T A A \mathbf{v} + \mathbf{v}^T A \mathbf{v} - \text{tr} A \rightarrow \max_A; \\ A^{-1} - I \succeq 0; \\ A^T = A. \end{cases} \quad (8)$$

Используя условие симметричности матрицы A , представим ее в виде $A = Q \Lambda Q^T$, где $Q^T = Q^{-1}$,

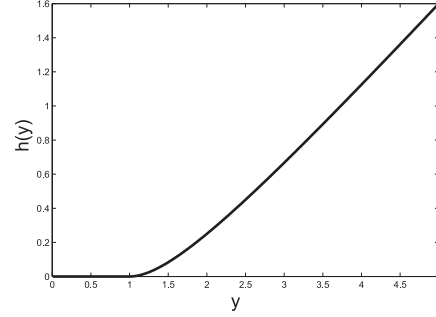


Рис. 1.

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ и $(\lambda_1, \dots, \lambda_m)$ — набор собственных чисел матрицы A . Заметим, что такое разложение единственно. Обозначим $\mathbf{x} = Q^T \mathbf{v}$. Тогда задача (8) переформулируется следующим образом:

$$\begin{cases} g(\mathbf{x}, \boldsymbol{\lambda}) = \sum_{j=1}^m (-\frac{1}{2} x_j^2 \lambda_j^2 + x_j^2 \lambda_j - \lambda_j) \rightarrow \max_{\mathbf{x}, \boldsymbol{\lambda}}; \\ \sum_{j=1}^m x_j^2 = \|\mathbf{v}\|^2; \\ 0 \leq \lambda_j \leq 1, \quad j = 1, \dots, m. \end{cases} \quad (9)$$

Максимальные значения функции $g(\mathbf{x}, \boldsymbol{\lambda})$ при фиксированных \mathbf{x} достигаются при следующих значениях $\boldsymbol{\lambda}$:

$$\lambda_j^*(\mathbf{x}) = \max\left(0, 1 - \frac{1}{x_j^2}\right). \quad (10)$$

Введем функцию $h(y)$ (см. рис. 1):

$$h(y) = \begin{cases} \frac{y}{2} + \frac{1}{2y} - 1, & y \geq 1; \\ 0, & 0 \leq y \leq 1. \end{cases} \quad (11)$$

Заметим, что функция $h(y)$ выпукла. Обозначим $y_j = x_j^2$. Тогда, подставив выражение (10) в систему (9), получим следующую задачу распределения ресурсов:

$$\begin{cases} \sum_{j=1}^m h(y_j) \rightarrow \max_{\mathbf{y}}; \\ \sum_{j=1}^m y_j = \|\mathbf{v}\|^2. \end{cases} \quad (12)$$

В случае $\|\mathbf{v}\| \leq 1$ критерий $\sum_{j=1}^m h(y_j)$ тождественно равен 0 (см. рис. 1), и компоненты y_j принимают любые значения от 0 до 1 при условии, что $\sum_{j=1}^m y_j = \|\mathbf{v}\|^2$. Таким образом, все $x_j^2 < 1$ и, следовательно, все $\lambda_j^* = 0$, то есть $A = 0$, и матрица регуляризации $R^{-1} = 0$. В этом случае регуляризатор штрафует любую попытку настройки на данные, и решающее правило становится вырожденным с $\mathbf{w}_{\text{MP}} = \mathbf{0}$. На практике такой случай соответствует ситуации, когда выбранное семейство

базисных функций $\{\varphi_i(\mathbf{x})\}_{i=1}^m$ не позволяет восстановить зависимость скрытой переменной от признаков. Заметим, что данный результат остается справедливым в подсемействах моделей, отвечающих диагональным и скалярным (пропорциональным единичной) матрицам регуляризации, т. к. вырожденная матрица регуляризации $R^{-1} = 0$ входит в эти подсемейства. Процесс поиска наилучшей модели в этих подсемействах является итеративным, поэтому возможность отсеять заведомо неадекватное подсемейство по аналитически проверяемому условию

$$\mathbf{w}_{ML}^T H \mathbf{w}_{ML} \leq 1 \quad (13)$$

позволяет значительно сократить время поиска наилучшей модели.

Если $\|\mathbf{v}\| > 1$, тогда, вследствие строгой выпуклости $h(y)$, максимум достигается в том случае, когда все ресурсы сосредоточены в одной из компонент y_i . Без ограничения общности выберем в качестве такой компоненты y_1 , т. е. $y_1 = \|\mathbf{v}\|^2$, $y_2 = \dots = y_m = 0$. Тогда первый собственный вектор матрицы A сонаправлен вектору \mathbf{v} . Остальные собственные векторы имеют нулевые собственные значения и могут быть выбраны произвольно при условии сохранения тождества $Q^T = Q^{-1}$. Отсюда, используя определение матрицы A , получим выражение для матрицы R^{-1} :

$$\begin{aligned} R^{-1} &= H^{-\frac{1}{2}}(A^{-1} - I)^{-1}H^{-\frac{1}{2}} = \\ &= H^{-\frac{1}{2}}Q \operatorname{diag}(\|\mathbf{v}\|^2 - 1, 0, \dots, 0)Q^T H^{-\frac{1}{2}} = \\ &= \frac{\|\mathbf{v}\|^2 - 1}{\|\mathbf{v}\|^2} H^{-\frac{1}{2}} \mathbf{v} \mathbf{v}^T H^{-\frac{1}{2}} = \\ &= \frac{\mathbf{w}_{ML}^T H \mathbf{w}_{ML} - 1}{\mathbf{w}_{ML}^T H \mathbf{w}_{ML}} \mathbf{w}_{ML} \mathbf{w}_{ML}^T. \end{aligned}$$

Таким образом, конечное выражение для матрицы R не зависит от выбора компоненты y_i при решении задачи распределения ресурсов (12), следовательно, исходная задача (7) имеет единственную точку максимума. Итак доказана следующая

Теорема 1. Функционал

$$f(R) = \log p(\mathbf{t} | X, \mathbf{w}_{MP}) - \operatorname{tr}(H(H+R)^{-1})$$

является унимодальным на множестве всех матриц $R = R^T \succeq 0$, если $\|H^{\frac{1}{2}} \mathbf{w}_{ML}\| > 1$.

Выводы

Доказанная теорема позволяет использовать точку максимума R критерия САИС для построения конечного решения задачи восстановления регрессии, гарантируя, что выбранная модель будет наилучшей среди всех допустимых моделей $R = R^T \succeq 0$.

Кроме того, полученный результат важен для выбора модели в задаче линейной регрессии, рассмотренной в [5]. В указанной работе подбор коэффициентов регуляризации, связанных индивидуально с каждым весом, производится путем максимизации непрерывного критерия Акаике (САИС). В этом случае критерий (5) оптимизируется на множестве всех диагональных матриц $R \succeq 0$. Унимодальность критерия на данном множестве остается открытой проблемой. Доказанная в данной работе теорема косвенно подтверждает предположение об унимодальности критерия и в этом случае. Полученное в работе условие релевантности (13) позволяет эффективно отсекают заведомо неадекватные наблюдаемым данным модели до начала итеративной настройки параметров модели. Аналитическое решение, максимизирующее непрерывный критерий Акаике, может быть использовано для построения решающих правил с лучшей обобщающей способностью.

Литература

- [1] Akaike H. A new look at statistical model identification // IEEE Trans. Automatic Control. 1974. V. 25. P. 461–464.
- [2] Kropotov D. A., Vetrov D. P. General solutions for information-based and bayesian approaches to model selection in linear regression and their equivalence // Pattern Recognition and Image Analysis. 2009. V. 3. P. 447–455.
- [3] Ezhova E., Mottl V., Krasotkina O. Estimation of time-varying linear regression with unknown time-volatility via continuous generalization of the Akaike information criterion // World Academy of Sciences, Engineering and Technology. 2009. V. 51.
- [4] Tipping M. E. The relevance vector machine // Advances Neural Information Processing Systems. 2000. V. 12. P. 652–658.
- [5] Kropotov D. A., Ptashko N. O., Vetrov D. P. Relevant regressors selection by continuous AIC // Pattern Recognition and Image Analysis. 2009. V. 3. Pp. 456–464.