

Сравнение эвристических алгоритмов выбора линейных регрессионных моделей*

Крымова Е. А., Стрижов В. В.

ekkrum@mail.ru, strijov@ccas.ru

Московский физико-технический институт,

Москва, Вычислительный центр РАН

В работе описан способ построения линейных регрессионных моделей, основанный на порождении и выборе признаков. Предложены эвристические алгоритмы выбора признаков. Выполнено сравнение этих алгоритмов с общеизвестными. Особенностью данного исследования является то, что задача выбора моделей поставлена для счетного набора признаков.

Введение

Процедура построения регрессионных моделей состоит из двух шагов. На первом шаге, на основе множества свободных переменных, порождается набор признаков. Один из способов построения такого набора описан в [1]. Модель-претендент есть линейная комбинация конечного подмножества признаков. На втором шаге производится выбор признаков, при этом выполняется настройка параметров модели и оценивается ее качество. Модель с настроенными параметрами, доставляющая минимум заданному функционалу качества, называется моделью оптимальной структуры.

Целью данной работы является сравнение предложенных эвристических алгоритмов с известными алгоритмами. Мотивацией работы является тот факт, что решение практических задач восстановления регрессионной зависимости требует рассмотрения большого числа порождаемых признаков. При таком условии алгоритмы, называемые «жадными», выбирают некоторый поднабор признаков, без возможности его модификации с целью улучшения структуры модели. Переборные алгоритмы, не обладая этим недостатком, имеют высокую вычислительную сложность. В данной работе предлагается компромиссный вариант алгоритма выбора регрессионных моделей и сравниваются следующие алгоритмы:

- 1) LARS метод наименьших углов [2];
- 2) полный перебор моделей [3];
- 3) метод группового учета аргументов [1];
- 4) алгоритм Лассо [4];
- 5) стохастическая структурная оптимизация;
- 6) шаговая регрессия [5, 6, 7];
- 7) оптимальное прорезивание в шаговой регрессии [8, 9];
- 8) модифицированный метод наименьших углов.

Сравнение алгоритмов показано на примере прикладной задачи, связанной с моделированием волатильности опционов по реальным историческим данным торгов опционом Brent Crude Oil.

*Работа выполнена при финансовой поддержке РФФИ, проекты № 07-07-00181, № 08-01-12022.

Постановка задачи

Задана выборка $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ — множество m пар, состоящих из вектора значений n свободных переменных $\mathbf{x}_i = (x_i^j)_{j=1}^n \in \mathbb{R}^n$, и значения одной зависимой переменной $y_i \in \mathbb{R}^1$. Индекс i элементов выборки и индекс j свободных переменных далее будем рассматривать как элементы множеств $i \in I = \{1, \dots, m\}$ и $j \in J = \{1, \dots, n\}$.

Задан класс регрессионных моделей $\mathcal{F} = \{f_s\}$ — параметрических функций, линейных относительно параметров,

$$y_i = f_s(\boldsymbol{\beta}_s, \mathbf{x}_i) = \sum_{j \in J_s} \beta_j x_i^j, \quad (1)$$

в которой $s \in \{1, \dots, 2^n\}$ является индексом модели, $\boldsymbol{\beta}_s = (\beta_j)_{j \in J_s}$ — вектор параметров, заданный индексом модели, $J_s \subseteq J$ — набор индексов свободных переменных. Введено ограничение на число элементов линейной комбинации (1). В множество \mathcal{F} могут входить только модели с числом свободных переменных $|J_s| \leq R$.

Принята следующая гипотеза порождения данных. Пусть случайная аддитивная переменная ν регрессионной модели

$$y = f(\boldsymbol{\beta}, \mathbf{x}) + \nu$$

имеет нормальное распределение $\mathcal{N}(0, \sigma_\nu^2)$.

Тогда, с учетом гомоскедастичности регрессионных остатков, распределение зависимой переменной имеет вид

$$p(y|x, \boldsymbol{\beta}, \sigma_\nu^2, f) = \frac{\exp(-\frac{1}{\sigma_\nu^2} S(D|\boldsymbol{\beta}, f))}{(2\pi\sigma_\nu^2)^{\frac{n}{2}}},$$

где S — сумма квадратов невязок $y_i - f(\boldsymbol{\beta}, \mathbf{x}_i)$. Это распределение задает указанный ниже критерий качества модели.

Дополнительно задано разбиение выборки $I = I^T \sqcup I^C$ на обучающую и контрольную. Для каждого набора данных, рассматриваемого в вычислительном эксперименте, наборы индексов I^T, I^C определены до начала эксперимента. Алгоритм выбора модели определяет метод оптимизации, доставляющий оптимальное значение параметрам $\tilde{\boldsymbol{\beta}}$

модели f на обучающей выборке $\{(x_i, y_i): i \in I^T\}$. Принят критерий качества — сумма квадратов регрессионных остатков на контрольной выборке

$$S = \sum_{i \in I^C} \left(y_i - f(\tilde{\beta}, x_i) \right)^2. \quad (2)$$

Требуется найти такую модель $f_s \in \mathcal{F}$, которая доставляет наименьшее значение функционалу качества. Такая модель будет называться моделью оптимальной структуры.

Порождение свободных переменных

Предлагается следующий способ формирования выборки D , состоящий из двух шагов.

Шаг первый. Задано множество непорождаемых свободных переменных $\Xi = \{\xi^u\}_{u=1}^U$. Задано конечное множество функций $G = \{g_v\}_{v=1}^V$. Рассмотрим декартово произведение $G \times \Xi$, элементу (g_v, ξ^u) которого поставлена в соответствие суперпозиция $g_v(\xi^u)$, однозначно определяемая индексами v, u . Обозначим $a_i = g_v(\xi^u)$, где индекс $i = (v-1)U + u$.

Шаг второй. Назначается базовая модель порождения признаков. В качестве модели, описывающей отношение между зависимой переменной y и свободными переменными a_i , используется полином Колмогорова-Габор:

$$y = \beta_0 + \sum_{i=1}^{UV} \beta_i a_i + \sum_{i=1}^{UV} \sum_{j=1}^{UV} \beta_{ij} a_i a_j + \dots,$$

где вектор коэффициентов

$$\beta = (\beta_0, \beta_i, \beta_{ij}, \beta_{ijk}, \dots)_{i,j,k,\dots=1,\dots,m}.$$

Запишем вышеприведенный ряд в виде

$$y = \sum_{j \in J} \beta_j x^j.$$

Переменные $\{x^j\}$ поставлены в однозначное соответствие мономам полинома.

Стандартизация данных

Выборка D стандартизирована таким образом, чтобы для $j \in J$ выполнялись условия нормировки

$$\sum_{i=1}^m x_i^j = 0, \quad \sum_{i=1}^m (x_i^j)^2 = 1, \quad \sum_{i=1}^m y_i = 0.$$

Предполагается, что векторы $\bar{x}_j = (x_1^j, \dots, x_m^j)$ и $\bar{x}_k = (x_1^k, \dots, x_m^k)$ для всех $j, k \in J$, $j \neq k$ линейно независимы.

LARS (метод наименьших углов)

LARS — алгоритм отбора признаков линейной модели [2]. На каждом шаге алгоритма происходит

изменение вектора параметров модели так, чтобы доставить добавляемому признаку наибольшую корреляцию с вектором регрессионных остатков. Основным достоинством LARS является то, что он выполняется за число шагов, равное числу свободных переменных.

Лассо

Лассо — алгоритм оценивания коэффициентов линейной модели [4]. Введение ограничения на сумму абсолютных значений коэффициентов модели приводит к обращению в 0 некоторых коэффициентов модели. Ненулевые коэффициенты соответствуют признакам, входящим в модель.

Обозначим сумму модулей коэффициентов модели $T(\beta) = \sum_{j=1}^n |\beta_j|$. Вектор коэффициентов $\hat{\beta}$ есть решения задачи минимизации $S(\beta)$ при ограничении: $T(\beta) \leq t$, где t — параметр регуляризации. Для решения задачи используется метод квадратичного программирования.

Шаговая регрессия

Шаговая регрессия — алгоритм последовательного удаления/добавления признаков [5]. Алгоритм последовательного добавления признаков присоединяет к текущему набору A по одному признаку, который доставляет максимум нижеприведенному критерию,

$$\hat{j} = \arg \max_{j \in J} \frac{S(A) - S(A \cup \bar{x}_j)}{S(A \cup \bar{x}_j)}.$$

Начальным считается пустой набор признаков.

Алгоритм последовательного удаления признаков начинает с самого большого набора, состоящего из всех признаков. На каждом шаге происходит удаление признака так, чтобы значение нижеприведенного критерия было как можно меньше:

$$\hat{j} = \arg \min_{j \in J} \frac{S(A \setminus \bar{x}_j) - S(A)}{S(A)}.$$

Останов алгоритма производится по выполнению условия C_p [10].

Алгоритм полного перебора

Этот алгоритм порождает все возможные множества мономов $\{\bar{x}_j\}_{j \in J}$. Пусть сложность модели $|J| \leq R$. Под сложностью модели понимается число линейно входящих параметров. Алгоритм последовательно строит модели-претенденты неубывающей сложности. Параметры каждой модели настраиваются методом наименьших квадратов по обучающей выборке. Наилучшая модель выбирается исходя из минимума ошибки на контрольной выборке. Введем переменную выбора монома —

вектор $\mathbf{c} = (c_1, \dots, c_n)$. Его элемент $c_j \in \{0, 1\}$ принимает значение 1, если $j \in J_s$, в противном случае 0. Базовая модель данного алгоритма имеет вид

$$y = \sum_{j \in J_s} c_j \beta_j x^j. \text{ Сложность этого алгоритма } \sum_{i=1}^R C_n^i.$$

Метод группового учета аргументов

Алгоритмы МГУА воспроизводят схему массовой селекции [1, 3]: последовательно порождаются модели возрастающей сложности. Каждая модель настраивается методом наименьших квадратов. Остановка алгоритма происходит, когда с увеличением номера шага начинается увеличение ошибки на контрольный выборке.

Стохастическая структурная оптимизация

Предложенный эвристический алгоритм состоит из итеративно выполняемых шагов. На первом шаге из множества признаков выбирается заданное число поднаборов, доставляющее соответствующей линейной модели наименьшее значение функционала качества. На втором шаге на заданном числе пар выполняется операция обмена признаками. На третьем шаге производится случайная замена произвольных признаков вновь полученных поднаборов. Шаги 2 и 3 итеративно повторяются. Алгоритм завершает работу, когда число шагов превысит заданное или когда ошибка оптимальной модели на контрольной выборке станет меньше заданной.

Оптимальное прореживание в шаговой регрессии

Оптимальное прореживание — это метод упрощения структуры регрессионной модели. Основная идея прореживания заключается в том, что те элементы модели, которые оказывают малое влияние на ошибку аппроксимации, можно исключить из модели без значительного ухудшения качества аппроксимации [8, 9].

Для построения регрессии требуется найти такие параметры β , которые доставляли бы наименьшее значение функции ошибки $S(\beta)$.

Локальная аппроксимация функции S в окрестности точки $\hat{\beta}$ с помощью разложения в ряд Тейлора записывается в виде

$$S(\beta + \Delta\beta) = S(\beta) + \mathbf{g}^T(\beta) \Delta\beta + \frac{1}{2} \Delta\beta^T H \Delta\beta + o(\|\beta\|^3),$$

где $\Delta\beta$ — возмущение вектора параметров β , $\mathbf{g} = \frac{\partial S}{\partial \beta}$ — градиент, $H = H(\beta) = \frac{\partial^2 S}{\partial \beta^2}$ — матрица вторых производных (матрица Гессе).

Функция $S(\beta)$ достигает своего максимума при $\beta = \hat{\beta}$, и ее поверхность квадратична. Таким образом, предыдущее выражение можно представить в виде $\Delta S = S(\beta + \Delta\beta) - S(\beta) = \frac{1}{2} \Delta\beta^T H \Delta\beta$.

Пусть исключение элемента модели есть исключение одного параметра модели, β_i . Исключенный параметр будем считать равным нулю. Исключение элемента эквивалентно выражению $\Delta\beta_i + \beta_i = 0$, иначе $\mathbf{e}_i^T \Delta\beta + \beta_i = 0$, где \mathbf{e}_i — вектор, i -й элемент которого равен единице, все остальные элементы равны нулю.

Требуется минимизировать квадратичную форму $\Delta\beta^T H \Delta\beta$ относительно $\Delta\beta$ при ограничениях $\mathbf{e}_i^T \Delta\beta + \beta_i = 0$ для всех значений i . Задача условной минимизации решается с помощью введения Лагранжиана $L = \Delta\beta^T H \Delta\beta - \lambda(\mathbf{e}_i^T + w_i)$, в котором λ — множитель Лагранжа. Дифференцируя Лагранжиан по приращению параметров и приравнявая его к нулю получаем (для каждого индекса i параметра β_i)

$$\Delta\beta = -\frac{\beta_i}{[H^{-1}]_{ii}} H^{-1} \mathbf{e}_i.$$

Этому значению вектора приращений параметров соответствует минимальное значение Лагранжиана

$$L_i = \frac{\beta_i^2}{2[H^{-1}]_{ii}}.$$

Полученное выражение называется мерой выпуклости функции ошибки S при изменении параметра β_i .

Функция L_i зависит от квадрата параметра β_i . Это говорит о том, что параметр с малым значением будет удален из модели. Однако если величина $[H^{-1}]_{ii}$ достаточно мала, это означает, что данный параметр оказывает существенное влияние на качество аппроксимации модели.

EM+LARS

В данной работе предлагается алгоритм, сочетающий в себе жадную стратегию LARS и перебор моделей. Это позволяет улучшить структуру модели, несущественно увеличив вычислительные затраты. Происходит порождение подмножеств признаков с помощью EM-алгоритма. На этих подмножествах с помощью LARS находятся параметры моделей, производится перебор полученных моделей. Модель, доставляющая наименьшую среднеквадратичную ошибку на контрольной выборке, считается оптимальной.

Задано K натуральных чисел в порядке возрастания C_k , $k = 1, \dots, K$, $C_k < n$. Пусть M — число моделей, получаемых на каждом шаге алгоритма. Рассмотрим k -й шаг алгоритма. Разобьем множество признаков X на C_k кластеров с помощью EM-алгоритма. M раз выбираем случайным образом из каждого кластера по одному элементу. Получаем M подмножеств из C_k признаков, принадлежащих разным классам. К каждому из подмножеств признаков на обучающей выборке приме-

Таблица 1. Результаты работы стандартных методов отбора признаков.

Алгоритм	$\frac{S(I^T)}{\ \mathbf{y}_T\ ^2}$	$\frac{S(I^C)}{\ \mathbf{y}_C\ ^2}$	AIC	k
LARS+EM	0.024	0.053	-863	9
Прорежив.	0.013	0.034	-1109	15
Стохаст.	0.014	0.024	-1299	20
Lasso	0.011	0.092	-491	11
Шаг. регр.	0.032	0.086	-405	30
МГУА	0.018	0.080	-584	8
LARS	0.019	0.089	-579	5
Перебор R=5	0.024	0.036	-	5

нием алгоритм LARS. В результате будут получены векторы параметров β_{kl} , $l = 1, \dots, M$.

За K шагов алгоритм находит KM моделей, выбирается оптимальная модель.

Вычислительный эксперимент

Сравнительный анализ алгоритмов был выполнен на исторических данных торгов опционом Brent Crude Oil [11]. Срок действия опциона — полгода, с 02.01.2001 по 26.06.2001. Тип опциона — put (право на продажу базового инструмента), символ CLG01. Базовый инструмент — нефть, символ NYM. Использовались ежедневные цены закрытия опциона и базового инструмента. Сетка цен исполнения опциона $\mathcal{K} = \{19.0, 19.5, \dots, 28.0, 28.5\}$.

Регрессионная выборка

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^m = \{(\langle K_i, t_i \rangle, \sigma_i)\}_{i=1}^m$$

была построена по исходным данным — историческим ценам опциона $C_{K,t}$ и базового инструмента P_t , где $K \in \mathcal{K}$, $t \in T$, следующим образом. Для каждого значения K_i и t_i , $i = 1, \dots, m$ вычислялось значение предполагаемой волатильности σ_i

$$\sigma_i = \arg \min_{\sigma \in [0, 1.5]} (C_{K_i, t_i} - C(\sigma, P_{t_i}, B, K_i, t_i)),$$

где справедливая цена опциона C определялась по формуле Блэка-Шоулза [12]. Длина истории составляла 112 отсчетов времени.

Множество порождающих функций было задано следующим образом: $G = \{1/x, \sqrt{x}, e^x\}$.

Максимальная степень полинома Колмогорова-Габора была выбрана равной трём. При этом число мономов составило 84. Регрессионная выборка была случайным образом разбита на контрольную и обучающую, равные по мощности. Стандартизация контрольной и обучающей выборок была проведена отдельно. Значения ошибок на обучении и контроле были усреднены по 10 запускам алгоритмов на различных разбиениях.

Результаты экспериментов представлены в Таблице 1. Для каждого алгоритма были вычисле-

ны: ошибки $S(I^T)$ и $S(I^C)$ на обучении и контроле (2), отнесенные к квадрату нормы соответствующего вектора ответов $\|\mathbf{y}_T\|^2 = \sum_{j \in I^T} y_j^2$ и $\|\mathbf{y}_C\|^2 = \sum_{j \in I^C} y_j^2$; значение информационного критерия

Акаике $AIC = m \ln \frac{S}{m} + 2k$; сложность модели k .

Заключение

В работе выполнено сравнение предложенных алгоритмов (стохастическая структурная оптимизация, модифицированный метод наименьших углов EM+LARS) с известными алгоритмами. Вычислительный эксперимент показал, что увеличение числа признаков позволяет добиться улучшения качества модели. Результаты экспериментов подтвердили жадность алгоритма LARS и большую эффективность алгоритма EM+LARS по сравнению с LARS. По результатам экспериментов наилучшими по совокупности критериев являются EM+LARS и алгоритм оптимального прореживания.

Литература

- [1] Malada H. R., Ivakhnenko A. G. Inductive Learning Algorithms for Complex Systems Modeling. CRC Press. 1994.
- [2] Efron B., Hastie T., Johnstone I., Tibshirani R. Least Angle Regression // The Annals of Statistics. 2004. Vol. 32, No. 2. Pp. 407–499.
- [3] Ивахненко А. Г., Степанко В. С. Помехоустойчивость моделирования. Киев: Наукова думка. 1985.
- [4] Tibshirani R. Regression shrinkage and Selection via the Lasso // Journal of the Royal Statistical Society. 1996. Vol. 32, No. 1. Pp. 267–288.
- [5] Draper N., Smith H. Applied Regression Analysis. John Wiley and Sons. 1981. Pp. 307–312.
- [6] Efrogmson M. A. Multiple regression analysis. Mathematical Methods for Digital Computers. Ralston, Wiley, New York. 1960.
- [7] Rawlings J. O., Pantula S. G., Dickey D. A. Applied Regression Analysis: A Research Tool. Springer-Verlag, New York. 1998.
- [8] Стрижов В. В. Поиск параметрической регрессионной модели в индуктивно заданном множестве // Журнал вычислительных технологий. 2007. № 1. С. 93–102.
- [9] Хайкин С. Нейронные сети, полный курс. М: Вильямс. 2008.
- [10] Mallows C. L. Some Comments on C_p . Technometrics. 1973. No. 15. Pp. 661–675.
- [11] Шуряев А. Н. Основы стохастической финансовой математики. Том 1. Факты. Модели. ФАЗИС. 2004.
- [12] Hull J. C. Options, Futures and Other Derivatives. Prentice Hall. 2000.