

АЛГОРИТМЫ ПОРОЖДЕНИЯ И ВЫБОРА РЕГРЕССИОННЫХ МОДЕЛЕЙ

Стрижов В.В., Крымова Е.А¹

Вычислительный центр РАН, e-mail: strijov@ccas.ru

¹Московский физико-технический институт, e-mail: ekkrym@mail.ru

Рассматриваются алгоритмы порождения и выбора линейных регрессионных моделей. Модель оптимальной структуры отыскивается как линейная комбинация элементов заданной базовой модели. Критерием оптимальности служит среднеквадратичная ошибка на тестовой подвыборке.

Цель работы — анализ и сравнение эвристических алгоритмов порождения моделей: однослойного и многослойного алгоритмов МГУА [1], метода стохастической структурной оптимизации и метода оптимального прореживания [2]. При анализе используются алгоритмы выбора свободных переменных: шаговая регрессия, метод наименьших углов (LARS) и лассо Тибширани.

Решается следующая задача. Задана выборка $\{(x_n, y_n)\}_{n=1}^N$, $\mathbf{x} \in \mathbf{R}^m$, $y \in \mathbf{R}$, которая разбита на обучающую и тестовую подвыборки случайным образом. Разбиение определено множествами индексов ℓ и C .

Принята базовая модель — полином Колмогорова-Габора, порождающая модели-претенденты: $y = \omega_0 + \sum_{i=1}^m \omega_i x_i + \sum_{i=1}^m \sum_{j=1}^m \omega_{ij} x_i x_j + \dots + \sum_{i=1}^m \dots \sum_{z=1}^m \omega_{i\dots z} \underbrace{x_i \dots x_z}_R$.

В этой модели $\mathbf{x} = \{x_i | i = 1, \dots, m\}$ — множество свободных переменных; $\boldsymbol{\omega}$ — вектор параметров $\boldsymbol{\omega} = \langle \omega_i, \omega_{ij}, \omega_{ijk}, \dots | i, j, k, \dots = 1, \dots, m \rangle$ и $F_0 = F_0(R)$ — число мономов.

Базовая модель представима в виде $\mathbf{y} = A\boldsymbol{\omega}$, где столбцы матрицы — значения мономов на выборке и $\mathbf{y} = \{y_1, \dots, y_N\}$. Обозначим соответствующие разбиения выборки как A_ℓ, y_ℓ и A_C, y_C .

Требуется выбрать такие столбцы матрицы A , задающие модель, которые доставляют минимум критерию оптимальности. Задача построения линейной регрессионной модели оптимальной структуры имеет вид $\mathbf{c} = \arg \min_{\mathbf{c} \in \{0,1\}^{F_0}} \|A_C(\boldsymbol{\omega} \times \mathbf{c}) - \mathbf{y}_C\|$, где \times — знак поэлементного умножения векторов. Параметры $\boldsymbol{\omega}$ определены как $\boldsymbol{\omega} = \arg \min_{\boldsymbol{\omega} \in \mathbf{R}^{F_0}} \|A_\ell \boldsymbol{\omega} - \mathbf{y}_\ell\|$.

На исторических данных опционных торгов был проведен анализ предложенных методов. Работа выполнена при поддержке РФФИ, проект № 07-07-00181.

Литература.

1. Malada H. R., Ivakhnenko A. G. Inductive Learning Algorithms for Complex Systems Modeling. CRC Press. 1994.
2. Стрижов В. В. Поиск параметрической регрессионной модели в индуктивно заданном множестве. Журнал вычислительных технологий. 2007. № 1. С. 93–102.