

А. Н. Тырсин, О. В. Ворфоломеева (Челябинск, ЮУрГУ).
Алгоритм непараметрического оценивания индекса детерминации регрессионных зависимостей.

Форма регрессионной модели во многих случаях не известна. Неправильный выбор формы может привести к ухудшению статистических характеристик уравнения, а также делает невозможным использование построенной модели за пределами выборки. Для оценки адекватности модели ее целесообразно сравнить с эмпирической регрессией. Это можно сделать с помощью сравнения индексов детерминации этих зависимостей.

Индекс детерминации показывает долю дисперсии результативной переменной Y , объясненной вариацией факторных переменных X_1, X_2, \dots, X_m , включенных в нелинейную модель регрессии

$$R_{Y/X_1 \dots X_m}^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad (1)$$

где \bar{y} — среднее значение результативной переменной, \hat{y}_i — значения регрессии, y_i — фактические значения переменной Y .

У эмпирической регрессии функциональная форма не известна, поэтому значения \hat{y}_i определяют без построения уравнения регрессии. В [1] рассмотрен ряд непараметрических алгоритмов построения регрессии. Они основаны на сглаживании значений y_i . Основным недостатком данного подхода является проблема выбора размера окрестности усреднения (апертуры сглаживающего фильтра). Описанный ниже алгоритм устраняет данный недостаток.

Пусть имеем многомерную выборку $(x_{i1}, \dots, x_{im}, y_i)$, $i = 1, \dots, n$. Формируем матрицу расстояний \mathbf{R} , элементы которой $r_{ij} = \sqrt{\sum_{k=1}^m (x_{ik} - x_{jk})^2}$, $i, j = 1, \dots, n$.

Для каждого i -го наблюдения подбираем оптимальную выборку из L_i ближайших соседей согласно матрице \mathbf{R} , так чтобы построенное уравнение линейной регрессии

$$\tilde{y}_l = a_{i0} + \sum_{k=1}^m a_{ik} x_{lk}, \quad l \in L_i = \{l_1, l_2, \dots, l_{L_i}\} \quad (2)$$

имело минимальную дисперсию ошибок регрессии

$$s_{L_i}^2 = \frac{\sum_{j \in L_i} (y_j - \tilde{y}_j)^2}{L_i - m - 1} \rightarrow \min_{3(m+1) < L_i \leq n}.$$

Для оптимальной локальной выборки из L_i ближайших соседей по формуле (2) формируем значение \hat{y}_i непараметрической регрессии

$$\hat{y}_i = a_{i0} + \sum_{k=1}^m a_{ik} x_{ik}.$$

По сформированному таким образом множеству значений \hat{y}_i , $i = 1, \dots, n$, по (1) найдем оценку индекса детерминации множественной регрессии.

Очевидно, что оптимальные размеры локальных выборок L_i в общем случае будут различными. Это позволит учесть изменение градиента теоретической функции регрессии в зависимости от значений факторных переменных и дисперсии случайной компоненты. Так как не использовалось никаких предположений о виде функции регрессии, то описанный алгоритм является непараметрическим.

Отметим, что оценку коэффициентов регрессии в (2) выполняют по-разному, в зависимости от особенностей исходных данных, например, можно использовать метод наименьших квадратов, метод наименьших модулей, а также робастные методы.

СПИСОК ЛИТЕРАТУРЫ

1. Хардле В. Прикладная непараметрическая регрессия. М.: Мир, 1993, 349 с.