

Кластерная регрессия

Дмитрий Федоряка

1 Введение

Здесь описывается алгоритм, предложенный автором в ходе написания статьи о композициях моделей в задаче векторной авторегрессии, но не показавший на практике улучшения точности в сравнении с базовым алгоритмом.

Предложенный алгоритм является алгоритмом композиции моделей. Композиция моделей — это алгоритм f , построенный на основе нескольких базовых алгоритмов f_1, \dots, f_K :

$$f(\mathbf{x}) = \sum_{i=1}^K w_i(\mathbf{x}) f_i(\mathbf{x})$$

Здесь f_1, \dots, f_K - базовые модели (в данной работе - линейные модели с различными матрицами весов). w_1, \dots, w_K - веса, с которыми модели учитываются. Для них должно выполняться $\sum_{i=1}^K w_i = 1$. Эти веса могут зависеть от объекта \mathbf{x} . В таком случае им можно придать вероятностный смысл: $w_i(\mathbf{x})$ - вероятность того, что объект \mathbf{x} описывается моделью f_i .

При построении алгоритма использовалась идея алгоритма кластеризации K-means: в пространстве объектов выделяются некоторые группы объектов и центры, затем за несколько итераций объекты распределяются между центрами так, чтобы каждый объект был ближе к «своему» центру, чем остальным. Для этого сначала строится произвольное разбиение объектов, а затем на каждой итерации центры вычисляются как центры масс в каждой группе и перераспределяются объекты.

Здесь вместо центров рассматриваются модели. Вместо вычисления центра масс производится обучение модели. «Ближайшую» модель для объектов находим как ту, которая даёт самую меньшую ошибку на этом объекте.

2 Описание алгоритма

2.1 Обучение

Опишем итерационный алгоритм, который строит модели f_i и находит веса w_i .

Произвольно разобьём все объекты на K непересекающихся подмножеств (классов): $1, \dots, m = \bigsqcup_{j=1}^K C_j$. Предположим, что объекты i -го класса описываются i -й моделью. Теперь запустим итерационный процесс, который должен осуществить это предположение.

На первом шагу итерации будем обучать модели на объектах соответствующих им классов.

На втором шагу будем вычислять для i -го объекта и j -й модели, вес $w(i, j)$, показывающий насколько хорошо данная модель описывает объект (в качестве критерия качества используя погрешность предсказания по отношению к известному истинному значению), и на основании полученных значений перераспределять объекты между классами.

После перераспределения объектов может возникнуть одна из двух исключительных ситуаций. В некоторый класс может не попасть ни один объект. Тогда, скорее всего, начальное количество классов K слишком большое, и надо удалить из рассмотрения эти

пустые классы, уменьшив K . Может оказаться, что перераспределения не произошло. Тогда нужно остановить алгоритм, т.к. на следующих итерациях перераспределения, скорее всего, происходить тоже не будет.

После некоторого числа итераций получим K разных моделей.

2.2 Предсказание

Пусть теперь нам надо найти ответ для объекта x . Для этого надо определить веса $w_i(x)$. Построим предварительный ответ для объекта, используя единственную линейную модель f_0 , обученную на всей выборке. Затем предскажем ответ с помощью каждой модели. Сравнивая предварительный ответ и ответы моделей, найдём веса. Для этого введём монотонно убывающую функцию $\mathfrak{E}(S_i(x))$, которая будет для ошибки $S_i(x) = \frac{\|f_i(x) - f_0(x)\|}{\|f_0(x)\|}$ вычислять меру правдоподобия. Вычислим правдоподобие для каждого класса, отнормируем - и получим веса моделей:

$$w_i(x) = \frac{\mathfrak{E}(S_i(x))}{\sum_{i=1}^K \mathfrak{E}(S_i(x))}$$

Примеры возможных функций $\mathfrak{E}(S)$: $\frac{1}{S+\varepsilon}$, e^{-S} , $1 - \frac{1}{1+e^{-S}}$.

3 Формальная запись

Algorithm 1 Кластерная регрессия

Вход:

$(X \in \mathbb{R}^{m \times n}, Y \in \mathbb{R}^{m \times r}, X^0 \in \mathbb{R}^n)$ — входные данные задачи авторегрессии;

K — число моделей;

N_{it} — число итераций;

\mathfrak{E} — функция преобразования погрешности в правдоподобие.

Выход: Y^0 — ответ задачи авторегрессии

- 1: $\forall j \in \overline{1, m} \ C_j = rand(1, m)$ — начальное разбиение на классы
 - 2: **для** $iter \in \overline{1, N_{it}}$
 - 3: **для** $j \in \overline{1, K}$ // М-шаг
 - 4: Обучаем алгоритм f_j на объектах с индексами C_j .
 - 5: **для** $i \in \overline{1, m}$ // Е-шаг
 - 6: **для** $j \in \overline{1, K}$
 - 7: $S(i, j) = \|f_j(X_i) - Y_i\|$ — вычисление ошибок
 - 8: **для** $i \in \overline{1, m}$ // Перераспределение
 - 9: **для** $j \in \overline{1, K}$
 - 10: $C_j = \{i \in \overline{1, m} | \underset{k}{argmin}(S(i, k)) = j\}$
 - 11: **если** $\exists i : C_i = \emptyset$ **то**
 - 12: $\{C_1, \dots, C_K\} = \{C_i | C_i \neq \emptyset\}$
 - 13: **если** перераспределение не поменяло класс ни для одного объекта **то**
 - 14: выйти из цикла
 - 15: Обучаем алгоритм f_0 на всей выборке.
 - 16: $\tilde{Y}^0 = f_0(X^0)$ — предварительный ответ.
 - 17: $\forall j \in \overline{1, K} w'_j = \mathfrak{E}(\frac{\|f_j(X^0) - \tilde{Y}^0\|}{\|\tilde{Y}^0\|})$ — вычисление вероятностей.
 - 18: $\forall j \in \overline{1, K} w_j = \frac{w'_j}{\sum_{k=1}^K w'_k}$ — нормировка.
 - 19: $Y^0 = \sum_{j=1}^K w_j \cdot f_j(X^0)$ — ответ.
-