

Об оптимизации билинейных форм в задачах классификации

В.В. Рязанов¹

¹Московский физико-технический институт (государственный университет)

В работе предлагается один алгоритм оптимизации билинейных форм в теории классификации. Подобные задачи возникают в ряде моделей, когда требуется выделение «информативной» подтаблицы таблицы обучения, то есть одновременное выделение в таблице обучения информативных признаков и объектов. Так в моделях вычисления оценок в одном из способов вычисления оценок объектов за классы оценки вычисляются как линейные функции параметров $\mathbf{x} = (x_1, x_2, \dots, x_n)$ «веса признаков» и одновременно линейные от параметров, характеризующих «веса объектов» - $\mathbf{y} = (y_1, y_2, \dots, y_m)$. В качестве функционала, характеризующего точность распознавания контрольной выборки, можно использовать эвристический функционал $f(\mathbf{x}, \mathbf{y}, t)$ - разность неотрицательных нормированных сумм оценок объектов за свои классы и аналогичной суммы оценок объектов за чужие классы, помноженной на скалярный неотрицательный параметр t . Данный параметр нужен для установления оптимального соответствия в функционале $f(\mathbf{x}, \mathbf{y}, t)$ между оценками за «свой» класс и оценками за «чужие» классы.

В алгоритмах АВО [1, 2] оценка $\Gamma_j(\mathbf{z})$ распознаваемого объекта $\mathbf{z} = (z_1, z_2, \dots, z_n)$ к классу K_j выражается как $\Gamma_j(\mathbf{z}) = \frac{1}{|K_j|} \sum_{\mathbf{z}_v \in K_j} y_v \sum_{\Omega \in \Omega_A} (\sum_{i \in \Omega} x_i) B_{\Omega}(\mathbf{z}_v, \mathbf{z})$, где $|K_j|$ - число эталонных объектов \mathbf{z}_v в классе K_j , Ω - некоторое опорное множество (подмножество признаков) из совокупности опорных множеств Ω_A алгоритма АВО, $B_{\Omega}(\mathbf{z}_v, \mathbf{z})$ - близость объекта \mathbf{z} к объекту обучения \mathbf{z}_v по опорному множеству. Можно показать, что оценка $\Gamma_j(\mathbf{z}) = \frac{1}{|K_j|} \sum_{\mathbf{z}_i \in K_j} y_i (\sum_{t \in J(\mathbf{z}, \mathbf{z}_i)} x_t) C_{d(\mathbf{z}, \mathbf{z}_i)-1}^{k-1}$, где

$J(\mathbf{z}, \mathbf{z}_i) = \{v : |z_{iv} - z_v| \leq \varepsilon_v, v = 1, \dots, n\}$, $d(\mathbf{z}, \mathbf{z}_i) = |J(\mathbf{z}, \mathbf{z}_i)|$, $k, \varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ - некоторые параметры алгоритма распознавания. Здесь в качестве опорных множеств используются всевозможные подмножества из k признаков. Тогда

$f(\mathbf{x}, \mathbf{y}, t) = \sum_{\alpha=1}^l \sum_{\mathbf{z}'_{\alpha} \in K_{\alpha}} \Gamma_{\alpha}(\mathbf{z}'_{\alpha}) - t \sum_{\alpha=1}^l \sum_{\mathbf{z}'_{\alpha} \notin K_{\alpha}} \Gamma_{\alpha}(\mathbf{z}'_{\alpha})$ (где \mathbf{z}'_{α} - элементы контрольной выборки), что является билинейной формой от параметров $x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_m$.

Итак, рассматривается следующая оптимизационная задача

$$f(\mathbf{x}, \mathbf{y}, t) = \sum_{i=1}^n \sum_{j=1}^m c_{ij}(t) x_i y_j \rightarrow \max_{\mathbf{x} \in X, \mathbf{y} \in Y} \quad (1)$$

Для выбора областей X и Y используются два варианта: $1 \geq x_i \geq 0, i = 1, 2, \dots, n, 1 \geq y_j \geq 0, j = 1, 2, \dots, m$, или $x_i \in \{0, 1\}, i = 1, 2, \dots, n, y_j \in \{0, 1\}, j = 1, 2, \dots, m$.

В последнем случае выбор единичных значений параметров $x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_m$ означает выбор некоторой подтаблицы таблицы обучения.

Для решения задачи (1) для случая дискретной оптимизации разработан (и исследован на модельных данных и данных из репозитория [3]) приближенный локальный метод, имеющий линейную сложность при малых окрестностях поиска.

Исследована также модификация модели поиска оптимальных $t, \mathbf{x}, \mathbf{y}$. Представляет интерес задача поиска оптимального значения параметра t , при котором будет наилучшее распознавание. Был использован стандартный критерий качества распознавания (доля правильно распознанных объектов контрольной выборки при выбранной подтаблице обучения) $F(\mathbf{x}, \mathbf{y}, t)$ и создан метод его максимизации.

Работа выполнена при поддержке гранта РФФИ № 16-31-00443.

Литература

1. Журавлев Ю.И., Никифоров В.В. Алгоритмы распознавания, основанные на вычислении оценок // Кибернетика. 1971. №3. С. 1-11.
2. Журавлев Ю.И. Об алгебраическом подходе к решению задач распознавания или классификации // Проблемы кибернетики. Вып. 33. М.: Наука, 1978. 5-68 с.
3. Lichman, M. UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>] // Irvine, CA: University of California, School of Information and Computer Science. 2013.