

УДК 519.584

М.П. Кузнецов, В.В. Стрижов, М.М. Медведникова

АЛГОРИТМ МНОГОКЛАССОВОЙ КЛАССИФИКАЦИИ ОБЪЕКТОВ, ОПИСАННЫХ В РАНГОВЫХ ШКАЛАХ

Рассматривается задача построения интегрального индикатора в ранговых шкалах [1,2,3]. В качестве практического приложения рассматривается проблема определения статуса угрожаемых видов животных, входящих в список Красной книги РФ [4]. В Красной книге РФ принята следующая категоризация редкости видов (таксонов) по степени угрозы их исчезновения. Имеется шесть различных категорий статуса (меток классов) таксонов: вероятно исчезнувшие, находящиеся под угрозой исчезновения, сокращающиеся в численности, редкие, неопределенные по статусу, восстанавливающиеся. Эта категоризация является монотонной: метки классов ранжированы по возрастанию биологического разнообразия. Назначение категории таксона может быть выполнено одним из методов согласования экспертных оценок [7,8], или путем аналитического вычисления категории на основе его описания, с учетом предложенной аналитиками модели [9].

Каждый таксон описан набором признаков, отражающих его состояние. Эксперт, владеющий информацией о таксоне, выставляет оценку для каждого признака в ранговой шкале. Таким образом, задана матрица «объект-признак», состоящая из описаний таксонов и вектор меток классов таксонов. Требуется построить модель, восстанавливающую класс таксона из Красной книги РФ по его описанию.

Задача ревизии Красной книги РФ и построения модели вычисления интегрального индикатора является актуальной из-за постоянного пополнения книги новыми записями о таксонах. Ранее были предложен ряд алгоритмов решения данной задачи [3,7]. Так как эти решения не включают процедуру выбора наиболее информативных признаков, ниже предлагается

поставить и решить задачу нелинейной коррекции экспертных оценок.

Постановка задачи

Задана множество $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$, $i \in \mathcal{I} = \{1, \dots, m\}$ пар. Каждая пара состоит из описания объекта \mathbf{x}_i (таксона) и соответствующей ему метки класса y_i (категория статуса таксона).

Описание объекта $\mathbf{x} = [\chi_1, \dots, \chi_j, \dots, \chi_n]$, $j \in \mathcal{J} = \{1, \dots, n\}$ — это набор экспертных оценок признаков. Оценки объектов по признакам выставлены в ранговых шкалах. Каждый признак i_j имеет собственную ранговую шкалу \mathbb{L}_j , состоящую из k_j упорядоченных элементов $\mathbb{L}_j = \{1 < 2 < \dots < k_j\}$. Значение класса y также принадлежит упорядоченному множеству $\mathbb{L}_0 = \{1 < 2 < \dots < k_0\}$.

Рассмотрим постановку задачи многоклассовой классификации в ранговых шкалах, включающую криволинейную модель $f(\mathbf{w}, \mathbf{x}_i)$ и соответствующую ей вектор-функцию $\mathbf{f}(\mathbf{w}, X) = [f(\mathbf{w}, \mathbf{x}_1), \dots, f(\mathbf{w}, \mathbf{x}_m)]$ с матрицей описаний $X = [\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_m]^T$ и зависимой переменной $\mathbf{y} = [y_1, \dots, y_i, \dots, y_m]$, где $\mathbf{w} = [w_1, \dots, w_s]$ — параметры модели. Эта модель должна доставлять минимум заданной функции ошибки $S(\mathbf{f}(\mathbf{w}, X), \mathbf{y})$.

Криволинейная модель $f(\mathbf{w}, \mathbf{x}_i)$ имеет вид

$$f(\mathbf{w}, \mathbf{x}_i) = \xi(\mathbf{b}_0, h(\mathbf{w}, \mathbf{x}_i)), \quad (1)$$

$$h(\mathbf{w}, \mathbf{x}_i) = \sum_{j \in \mathcal{J}} u_j g(\mathbf{b}_j, x_{ij}). \quad (2)$$

где вектор параметров $\mathbf{w} = [\mathbf{b}_0; \mathbf{b}_1; \dots; \mathbf{b}_n; \mathbf{u}] = [\mathbf{b}_0^T, \mathbf{b}_1^T, \dots, \mathbf{b}_n^T, \mathbf{u}^T]^T$ состоит из векторов \mathbf{b}_j — параметров монотонной коррекции j -го признака χ_j и весовых коэффициентов признаков $\mathbf{u} = [u_1, \dots, u_j, \dots, u_n]^T$. Функция g монотонной коррекции задана следующим образом:

$$g(\mathbf{b}_j, \chi): \chi \mapsto \mathbf{b}_j = \begin{cases} 1 \mapsto b_{j1}, \\ 2 \mapsto b_{j2}, \\ \dots \\ k_j \mapsto b_{jk_j}. \end{cases}$$

При этом соблюдается условие монотонности параметров,

$$\text{Ord}(\mathbf{b}_j): 0 < b_{j1} < b_{j2} < \dots < b_{jk_j} < 1 \quad \text{для } j = 1, \dots, n \quad \text{и} \quad \text{Ord}(\mathbf{b}_0): b_{01} < b_{02} < \dots < b_{0k_0}. \quad (3)$$

Функция $\xi(\mathbf{b}_0, h(\mathbf{w}, \mathbf{x}_i))$ определяет для числа $h(\mathbf{w}, \mathbf{x}_i)$ ближайшую по модулю компоненту вектора \mathbf{b}_0 :

$$\xi(\mathbf{b}_0, h(\mathbf{w}, \mathbf{x}_i)) = \arg \min_{j \in J} |b_{0j} - h(\mathbf{w}, \mathbf{x}_i)|.$$

Введя обозначение для матрицы скорректированных экспертных оценок

$$G = [g_{ij}] = [g(\mathbf{b}_j, x_{ij})], \quad i \in \mathcal{I}, j \in \mathcal{J},$$

перепишем (1) и (2) в виде модели интегрального индикатора

$$f(\mathbf{w}, \mathbf{x}_i) = \xi(\mathbf{b}_0, [G\mathbf{u}]_i). \quad (4)$$

Назначим функцией ошибки модели сумму квадратов регрессионных остатков, $S(\mathbf{w}) = \|\mathbf{f}(\hat{\mathbf{w}}, X) - \mathbf{y}\|_2^2 + \lambda \|\hat{\mathbf{u}}\|_2^2$, включающую регуляризующее слагаемое с фиксированным коэффициентом λ , где $\hat{\mathbf{w}}$ и $\hat{\mathbf{u}}$ — параметры, которые необходимо оценить.

Оценивание параметров модели

Оценивание параметров \mathbf{w} модели \mathbf{f} выполняется итеративно. Перед началом итераций значения векторов $\mathbf{b}_0, \mathbf{b}_1, \dots, \mathbf{b}_n$ назначены таким образом, что функция g является тождественной, $g = \text{id}$. Оценивание параметров выполняется в три шага. Сначала при фиксированных значениях векторов $\hat{\mathbf{b}}_0, \dots, \hat{\mathbf{b}}_n$ оцениваются весовые коэффициенты

$$\hat{\mathbf{u}} =_{\mathbf{u} \in \mathbb{R}^n} S([\hat{\mathbf{b}}_0; \dots; \hat{\mathbf{b}}_n; \mathbf{u}]).$$

Затем при фиксированных значениях коэффициентов $\hat{\mathbf{u}}$ оцениваются параметры монотонной коррекции

$$[\mathbf{b}_1; \dots; \mathbf{b}_n] =_{\text{Ord}(\mathbf{b}_1), \dots, \text{Ord}(\mathbf{b}_n)} S([\hat{\mathbf{b}}_0; \dots; \hat{\mathbf{b}}_n; \mathbf{u}])$$

с учетом требования монотонности (3) значений этих параметров. На последнем этапе оценивается вектор \mathbf{b}_0 : $\mathbf{b}_0 =_{\text{Ord}(\mathbf{b}_0)} S([\hat{\mathbf{b}}_0; \dots; \hat{\mathbf{b}}_n; \mathbf{u}])$.

Итерации выполняются до стабилизации функции ошибки S .

Рассмотрим эти три этапа более подробно. За начальное приближение примем столбцы матрицы G

$$\hat{G} = [\mathbf{g}(\hat{\mathbf{b}}_1, \chi_1), \dots, \mathbf{g}(\hat{\mathbf{b}}_n, \chi_n)] = [\chi_1, \dots, \chi_n],$$

поскольку, как было сказано выше, $g = \text{id}$, и вектор $\hat{\mathbf{y}} = \mathbf{y}$. Таким образом, векторы $\hat{\mathbf{b}}_0, \hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_n$ в начальном приближении в качестве элементов содержат элементы множеств $\mathbb{L}_0, \mathbb{L}_1, ts, \mathbb{L}_n$.

Шаг 1. Найдем $\hat{\mathbf{u}}$ при фиксированных значениях $\hat{\mathbf{b}}_0, \dots, \hat{\mathbf{b}}_n$:

$$\hat{\mathbf{u}} = \arg \min_u \|\hat{\mathbf{y}} - \hat{G}\mathbf{u}\| + \lambda \|\mathbf{u}\|.$$

Решение на шаге 1 имеет вид: $\hat{\mathbf{u}} = (\hat{G}^T \hat{G} + \lambda I)^{-1} \hat{G}^T \hat{\mathbf{y}}$.

Шаг 2. При фиксированных $\hat{\mathbf{b}}_0, \hat{\mathbf{u}}$ оценим скорректированную матрицу описаний $\mathbf{G} = [\mathbf{g}(\mathbf{b}_1, \chi_1), \dots, \mathbf{g}(\mathbf{b}_n, \chi_n)] = [\mathbf{g}_1, \dots, \mathbf{g}_n]$.

Для каждого $\mathbf{g}_j \in \mathbb{R}^m$ будем вычислять вектор $\hat{\mathbf{g}}_j$, являющийся монотонной коррекцией исходного вектора \mathbf{g}_j :

$$\begin{cases} [\hat{\mathbf{g}}_1, \dots, \hat{\mathbf{g}}_n] = \arg \min \|\xi(\mathbf{b}_0, G\hat{\mathbf{u}}) - \hat{\mathbf{y}}\|_2^2, \\ \text{из } g_{ij_1} \leq g_{ij_2} \text{ следует } \hat{g}_{ij_1} \leq \hat{g}_{ij_2} \quad i \in \mathcal{J}, j_1, j_2 \in \mathcal{J}, \\ g_{ij} \in [0, 1] \quad i \in \mathcal{J}, j \in \mathcal{J}, \text{ согласно (3)}. \end{cases}$$

По векторам $\hat{\mathbf{g}}_1, \dots, \hat{\mathbf{g}}_n$ затем однозначно восстанавливаются векторы $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_n$ как упорядоченные векторы, содержащие различные элементы $\hat{\mathbf{g}}_1, \dots, \hat{\mathbf{g}}_n$. Для решения этой задачи используется алгоритм градиентного спуска.

Шаг 3. При фиксированных $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_n, \hat{\mathbf{u}}$ оценим вектор \mathbf{b}_0 и $\hat{\mathbf{y}} = \mathbf{g}(\mathbf{b}_0, \mathbf{y})$:

$$\hat{\mathbf{b}}_0 = \arg \min_{\text{Ord}(\mathbf{b}_0)} \|\xi(\mathbf{b}_0, \hat{G}\hat{\mathbf{u}}) - \mathbf{g}(\mathbf{b}_0, \mathbf{y})\|_2^2.$$

Выбор признаков при классификации

Так как число объектов в данной задаче, определенное составом Красной книги РФ, сопоставимо с числом признаков, необходимо выбрать наиболее информативные признаки. Множество индексов признаков, включенных в модель, назовем активным набором и обозначим $\mathcal{A} \subseteq \mathcal{J}$.

Поставим задачу выбора наиболее информативных признаков следующим

образом. Разобьем выборку \mathcal{D} на две подвыборки, обучающую и тестовую. Обозначим индексы элементов этих подвыборок соответственно $\mathcal{L} \sqcup \mathcal{T} = \mathcal{J}$. Для некоторого активного набора признаков \mathcal{A} найдем на обучающей подвыборке $\mathcal{D}_{\mathcal{L}}$ оптимальные, согласно заданной функции ошибки S , параметры $\hat{\mathbf{w}}_{\mathcal{A}}$, $\hat{\mathbf{w}}_{\mathcal{A}} =_{\mathbf{w}} S(\mathbf{w}_{\mathcal{A}}|\mathcal{D}_{\mathcal{L}})$.

Затем выберем наиболее информативные признаки — активный набор $\hat{\mathcal{A}}$ по всем поднаборам индексов признаков $\mathcal{A} \subseteq \mathcal{J}$, доставляющий на тестовой выборке $\mathcal{D}_{\mathcal{L}}$ минимум функции ошибки: $\hat{\mathcal{A}} =_{\mathcal{A} \subseteq \mathcal{J}} S(\hat{\mathbf{w}}_{\mathcal{A}}|\mathcal{D}_{\mathcal{T}})$.

Для выбора наиболее информативного подмножества признаков используется итеративный алгоритм добавления признаков.

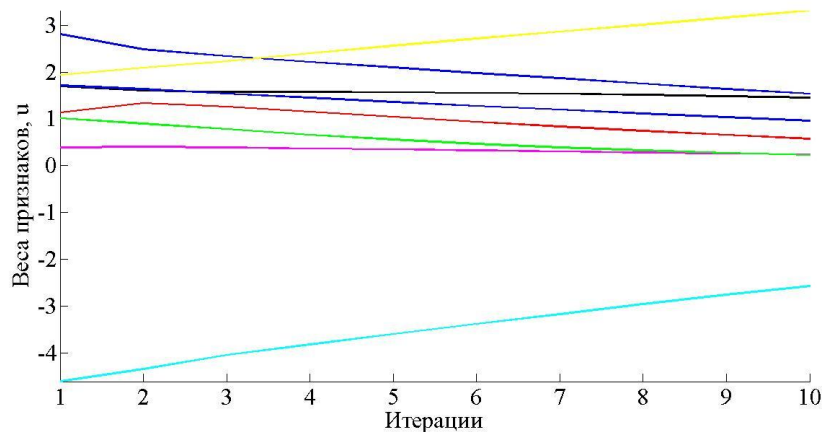


Рис. 1. Изменение весов признаков

На первом шаге этого алгоритма принимается активное множество информативных признаков $\hat{\mathcal{A}} = \emptyset$. На каждом следующем шаге к множеству $\hat{\mathcal{A}}$ добавляется признак с индексом \hat{j} , такой что $\hat{j} =_{j \in \mathcal{J} \setminus \mathcal{A}} S(\hat{\mathbf{w}}_{\hat{\mathcal{A}} \cup \{j\}}|\mathcal{D}_{\mathcal{T}})$.

Эта процедура продолжается итеративно до тех пор, пока значение функции ошибки S на контрольной выборке $\mathcal{D}_{\mathcal{T}}$ не достигнет минимума.

Вычислительный эксперимент

Работа алгоритма иллюстрируется данными из Красной Книги РФ. Экспертами заполнена таблица данных для 29 различных объектов. Каждый объект описывается 102 признаками. Отобрано восемь наиболее

информативных признаков. В качестве функции ошибки классификации принимается величина

$$Q = \frac{1}{m} \sum_{i=1}^m S(f(\mathbf{w}_{\mathcal{D} \setminus i}, X_{\mathcal{J} \setminus i}), y_i),$$

где $\mathbf{w}_{\mathcal{J} \setminus i}$ — вектор параметров, оцененный по всей выборке без учета i -го объекта, а $X_{\mathcal{J} \setminus i}$ — матрица X , в которой исключена i -я строка. Эта метрика представляет собой среднюю ошибку классификации на всех объектах выборки, ее значение составило $Q = 0,75$.

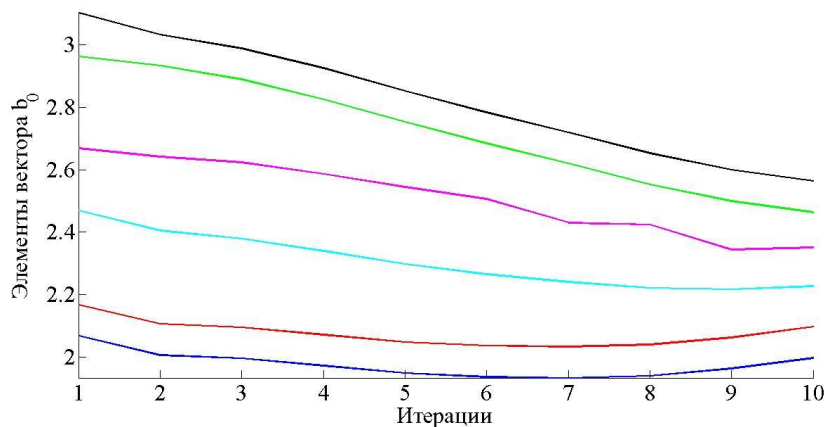


Рис. 2. Изменение элементов вектора параметров b_0 .

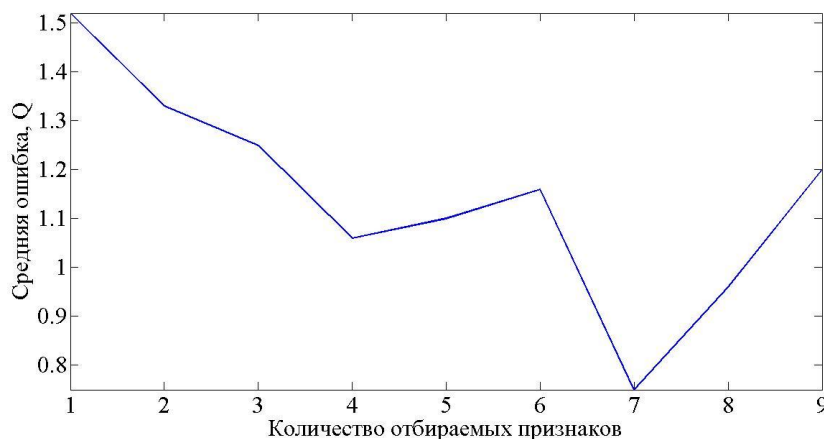


Рис. 3. Зависимость функции ошибки от количества выбираемых признаков.

Ниже представлены графики изменения параметров относительно итераций алгоритма. На рис. 1 показана изменение весов регрессии \mathbf{u} , на рис. 2 —

изменение элементов вектора b_0 . По оси абсцисс отложено количество итераций, по оси ординат — количественное значение каждого признака. Прекращение изменений наблюдается на десятой итерации.

На рис. 3 показана зависимость функции ошибки от количества выбираемых признаков. Видно, что ее минимум достигается при семи признаках, и значение средней ошибки равно $Q = 0.75$.

Заключение

Предложен метод построения рангового интегрального индикатора на примере задачи категоризации таксонов Красной книги РФ. Данный метод отличается от обычной задачи восстановления регрессии тем, что исходные данные представлены в ранговых шкалах и корректируются в процессе вычисления интегрального индикатора. Предложен алгоритм отбора наиболее информативных признаков. С помощью этого алгоритма получена адекватная модель получения категорий новых таксонов.

Список литературы

- [1] **Стрижов, В.В.** Уточнение экспертных оценок с помощью измеряемых данных // Заводская лаборатория. Диагностика материалов. – 2006. – Т. 72(7). – С. 59–64.
- [2] **Strijov, V.V. et al.** Integral indicator of ecological impact of the Croatian thermal power plants // Energy. –2011. – Vol. 36(7) . – Pp. 4144–4149.
- [3] **Стрижов, В. В.** Уточнение экспертных оценок, выставленных в ранговых шкалах, с помощью измеряемых данных // Заводская лаборатория. Диагностика материалов. –2011. –Т. 77(7). – С. 72–78.
- [4] Красная книга Российской Федерации (животные). М: АСТ Астрель, 2001.
- [5] **Литвак, Б.Г.** Экспертная информация: Методы получения и анализа. М.: Радио и связь. –1982.
- [6] **Орлов, А.И.** Организационно-экономическое моделирование. Экспертные оценки. М: МГТУ им. Н. Э. Баумана. – 2011.

[7] **Kotowski, W.** Rule learning with monotonicity constraints // Proceedings of the 26th Annual International Conference on Machine Learning. – 2009. – Vol 382. – P. 68.