

## Лекция 5

Прицип частичной прецедентности, тестовый алгоритм,  
модель АВО

*Лектор – Сенько Олег Валентинович*

Курс «Математические основы теории прогнозирования»  
4-й курс, III поток

- 1 Тестовый алгоритм
- 2 Представительные наборы
- 3 Алгоритмы вычисления оценок

Существует ряд методов распознавания, основанных на **принципе частичной прецедентности**. Данный принцип подразумевает поиск по обучающей выборке фрагментов описаний, позволяющих разделить распознаваемые классы  $K_1, \dots, K_L$ . Распознаваемый объект оценивается по совокупности найденных фрагментов. Одной из первых реализаций принципа частичной прецедентности является тестовый алгоритм, предложенный в 1966 году. Данный алгоритм основан на понятии тупикового теста. Исходный вариант тестового алгоритма предназначен для распознавания объектов, описываемых с помощью бинарных или категориальных признаков  $X_1, \dots, X_n$ . Иными словами  $X_i \in \{a_1^i, \dots, a_{k_i}^i\}$ ,  $i = 1, 2, \dots, n$ .

Пусть обучающая выборка  $\tilde{S}_t$ , содержит объекты из классов  $K_1, \dots, K_L$ . При этом общее число объектов равно  $m$ . Выборке  $\tilde{S}_t$  ставится в соответствие таблица  $\mathbb{T}_{nmL}$ ,  $j$ -ой строкой которой являются значения признаков для объекта  $s_j$ .

**Определение 1.** Тестом таблицы  $\mathbb{T}_{nmL}$  называется совокупность столбцов  $i_1, \dots, i_r$  таких, что после удаления из  $\mathbb{T}_{nmL}$  всех столбцов, за исключением  $i_1, \dots, i_r$ , в полученной таблице все пары строк, соответствующие разным классам различны хотя бы по одному признаку

**Определение 2.** Тест  $T = \{i_1, \dots, i_r\}$  называется тупиковым, если никакое его отличное от  $T$  подмножество (собственное подмножество) тестом не является.

На этапе обучения ищется множество всевозможных тупиковых тестов  $\tilde{T}(\tilde{S}_t)$  для таблицы  $\mathbb{T}_{nmL}$ . Предположим что нам требуется распознать объект  $s_*$  с векторным описанием  $(x_{*1}, \dots, x_{*n})$ .

Выделим в векторном описании фрагмент  $(x_{*i_1}, \dots, x_{*i_r})$ , соответствующий тесту

$$T = \{i_1, \dots, i_r\}, T \in \tilde{T}(\tilde{S}_t).$$

Фрагмент  $(x_{*i_1}, \dots, x_{*i_r})$  сравнивается с множеством фрагментов строк таблицы  $\mathbb{T}_{nmL}$ , соответствующих классу  $K_l$ :  $\{(x_{ji_1}^t, \dots, x_{ji_r}^t) \mid s_j \in K_l\}$ . В случаях, когда выполняются равенства

$$x_{*i_1} = x_{ji_1}^t, \dots, x_{*i_r} = x_{ji_r}^t$$

фиксируем полное совпадение распознаваемого объекта  $s_*$  с объектом  $s_j \in K_l$  тесту  $T$ . Обозначим число полных совпадений через  $G_l(\mathbb{T}_{nmL}, s_*)$ . Оценка объекта  $s_*$  за класс  $K_l$  вычисляется по формуле:

$$\gamma(s_*) = \frac{1}{m_l} \sum_{T \in \tilde{T}(\tilde{S}_t)} G_l(\mathbb{T}_{nmL}, s_*),$$

где  $m_l$  - число объектов обучающей выборки из класса  $K_l$ .

Классификация объекта  $s_*$  может производиться по вектору оценок  $[\gamma_1(s_*), \dots, \gamma_L(s_*)]$  с помощью стандартного решающего правила, т.е. объект относится в тот класс, оценка за который максимальна. Задача нахождения множества всех тупиковых тестов таблицы  $\mathbb{T}_{nmL}$  сводится к задаче поиска всех тупиковых покрытий матрицы сравнений  $\mathbb{C}_{nmL}$ , которая строится по матрице  $\mathbb{T}_{nmL}$ . Каждой паре классов  $K_l$  и  $K_{l'}$  в матрице  $\mathbb{C}_{nmL}$  сопоставлена подматрица  $\mathbb{C}_{nmL}^{ll'}$ , состоящая из  $m_l m_{l'}$  строк. Пусть  $(c_{f1}^{ll'}, \dots, c_{fn}^{ll'})$  строка матрицы  $\mathbb{C}_{nmL}^{ll'}$  соответствует сравнению описаний объекта  $x_g$  из класса  $K_l$  и описание объекта  $x_{g'}$  из класса  $K_{l'}$ . Элемент  $c_{fi}^{ll'} = 0$ , если  $x_{gi} = x_{g'i}$ , и  $c_{fi}^{ll'} = 1$ , если  $x_{gi} \neq x_{g'i}$ . Таким образом  $\mathbb{C}_{nmL}$  имеет размерность  $M \times n$ , где  $\sum_{l=1}^L \sum_{l'=1}^{l-1} m_l m_{l'}$ .

Мы будем говорить, что столбец с номером  $j$  матрицы  $C_{nmL}$  покрывает строку  $(c_{f1}, \dots, c_{fn})$ , если  $c_{fj} = 1$ . Набор столбцов  $\{j_1, \dots, j_r\}$  образует покрытие матрицы  $C_{nmL}$ , если при любом  $f \in \{1, \dots, M\}$  существует такое  $j' \in \{j_1, \dots, j_r\}$ , что  $c_{fj'} = 1$ . Покрытие  $\tilde{J}$  называется тупиковым, если его произвольное собственное подмножество, покрытием не является. Очевидно, что для произвольного набора столбцов обладание свойством тупикового набора для  $\mathbb{T}_{nmL}$  эквивалентно обладанию свойством тупикового покрытия для  $C_{nmL}$ . Таким образом задача о поиске всевозможных тупиковых тестов сводится к известной задаче о поиске всевозможных тупиковых покрытий. Нахождение всех тупиковых тестов является сложной комбинаторной задачей. Однако эффективные алгоритмы поиска разработаны для некоторых типов таблиц. При решении практических задач эффективен подход, основанный на вычислении только части тупиковых тестов.

Другим известным классом алгоритмов распознавания, основанным на принципе частичной прецедентности, являются алгоритмы типа КОРА. В отличие от тестового алгоритма, где в качестве информативных элементов используются несжимаемые наборы признаков – тупиковые тесты, в алгоритмах типа КОРА в качестве информативных элементов используются несжимаемые фрагменты описаний эталонных объектов обучающей выборки.

**Определение 3.** Пусть  $(x_{v1}, \dots, x_{vn})$  - признаковое описание объекта  $s_v \in \tilde{S}_t \cap K_l$ . Набор  $(x_{vj_1}, \dots, x_{vj_r})$  называется представительным набором для класса  $K_l$  при выполнении следующего условия. Пусть  $(x_{uj_1}, \dots, x_{uj_r})$  произвольная строка таблицы  $\mathbb{T}_{nmL}$ , которая соответствует объекту  $s_u$ , не принадлежащему  $\tilde{S}_t \cap K_l$ . Тогда существует такое  $j' \in \{j_1, \dots, j_r\}$ , что  $x_{uj'} \neq x_{vj'}$ .

**Определение 4.** Представительный набор называется **тупиковым**, если никакое его собственное подмножество представительным набором не является.

На этапе обучения для каждого из классов  $K_1, \dots, k_L$  по таблице  $T_{nmL}$  ищется множество всевозможных тупиковых представительных наборов. Пусть  $\tilde{V}_l$  - множество всевозможных представительных наборов для класса  $K_l$ .

Предположим, что нам требуется распознать объект  $s_*$  с описанием  $(x_{*1}, \dots, x_{*n})$ . Пусть  $\mathbf{v} = (x_{vj_1}, \dots, x_{vj_r})$  - представительный набором. Функция  $B(s_*, \mathbf{v})$  равна 1, при выполнении всех неравенств  $(x_{*j_1} = x_{vj_1}, \dots, x_{*j_r} = x_{vj_r})$  и  $B(s_*, \mathbf{v})$  равна 0 в противном случае. Оценка  $s_*$  за класс  $K_l$  вычисляется по формуле

$$\Gamma(s_*, K_l) = \frac{1}{|\tilde{V}_l|} \sum_{\mathbf{v} \in \tilde{V}_l} B(s_*, \mathbf{v})$$

Для нахождения тупиковых представительных наборов для класса  $K_l$ , содержащихся в эталонном описании  $x_v$  объекта  $s_v$  формируются матрица сравнения  $C_{nmL}^{lv}$  со всеми описаниями других классов таблицы  $T_{nmL}$ . Пусть строка  $(c_{f1}^{lv}, \dots, c_{fn}^{lv})$  матрицы  $C_{nmL}^{lv}$  соответствует сравнению  $x_v$  с описанием  $x_g$  объекта  $s_g$ , не принадлежащему  $K_l \cap \tilde{S}_t$ . Элемент  $c_{f1}^{lv} = 0$ , если  $x_{vj} = x_{gj}$ , и  $c_{f1}^{lv} = 1$ , если  $x_{vj} \neq x_{gj}$ . Таким образом матрица  $C_{nmL}^{lv}$  имеет размер  $(m - m_l)n$ . Тупиковые покрытия матриц сравнения  $C_{nmL}^{lv}$  определяют тупиковые представительные наборы, являющиеся фрагментами описания  $x_v$ . Полное множество представительных наборов для класса  $K_l$  является объединением множеств представительных наборов, найденных для описаний всех объектов обучающей выборки из  $K_l$ . Таким образом задача поиска всех представительных наборов сводится к решению  $m_l$  задач поиска тупиковых покрытий для матриц сравнения размера  $(m - m_l)n$ .

Первоначальные варианты тестового алгоритма и алгоритма типа КОРА были разработаны для бинарных или категориальных переменных. Они не могут быть напрямую использованы в задачах с признаками, принимающими значения из интервалов вещественной оси. Для того, чтобы обеспечить возможность работы с подобной информацией могут быть использованы два подхода.

Первый подход основан на разбиении области возможных значений каждого вещественнозначного признака на  $k$  связных подмножеств (интервалов, полуинтервалов, отрезков). Значению признака, принадлежащего  $j$ -ому элементу разбиения присваивается значение  $j$ . Разбиение оптимизируется с целью достижения максимального разделения классов. Выбирается такое число элементов разбиения  $k$ , при котором достигается максимальная точность распознавания.

Другой подход основан на модификации понятий теста и представительного набора с использованием пороговых параметров  $(\varepsilon_1, \dots, \varepsilon_n)$ , которые задаются для признаков  $X_1, \dots, X_n$ .

**Определение 5.** Тестом таблицы  $\mathbb{T}_{nmL}$  называется совокупность столбцов  $i_1, \dots, i_r$  таких, что после удаления из  $\mathbb{T}_{nmL}$  всех столбцов, за исключением столбцов  $i_1, \dots, i_r$  в полученной таблице  $\mathbb{T}_{rmL}$  для всех пар строк, соответствующих разным классам абсолютная величина различий хотя бы по одному признаку  $X_j$  превышает  $\varepsilon_j$ .

Аналогичным образом вводится модифицированное определение представительного набора.

**Определение 6.** Пусть  $(x_{v1}, \dots, x_{vn})$ - признаковое описание объекта  $s_v \in \tilde{S}_t \cap K_l$ . Набор  $(x_{vj_1}, \dots, x_{vj_r})$  называется представительным набором для класса  $K_l$  при выполнении следующего условия. Пусть  $(x_{uj_1}, \dots, x_{uj_r})$  произвольная строка таблицы  $\mathbb{T}_{nmL}$  соответствующая объекту  $s_u$ , не принадлежащему  $\tilde{S}_t \cap K_l$ . Тогда существует такое  $j' \in \{j_1, \dots, j_r\}$ , что  $|x_{uj'} - x_{vj'}| > \varepsilon_{j'}$ .

Главным требованием при выборе  $\varepsilon$ -порогов является достижение максимальной отделимости объектов разных классов при сохранении сходства внутри классов. Поиск тупиковых тестов и тупиковых представительных наборов при модифицированных определениях аналогичен их поиску в первоначальных вариантах методов.

Тестовый алгоритм и алгоритм с представительными наборами являются частью более общей конструкции, основанной на принципе частичной прецедентности и носящей название алгоритмов вычисления оценок. Существует много вариантов моделей данного типа. Причём конкретный вид модели определяется выбранными способами задания различных её элементов. Рассмотрим основные составляющие модели.

**Задание системы опорных множеств.** Под **опорными множествами** модели АВО понимается наборы признаков, по которым осуществляется сравнение распознаваемых и эталонных объектов. Примером системы опорных множеств является множество тупиковых тестов. Система опорных множеств  $\Omega_A$  некоторого алгоритма  $A$  может задаваться через систему подмножеств множества  $\{1, \dots, n\}$  или через систему характеристических бинарных векторов.

Каждому подмножеству  $\{1, \dots, n\}$  может быть сопоставлен бинарный вектор размерности  $n$ . Пусть  $\{i_1, \dots, i_k\} \subseteq \{1, \dots, n\}$ . Тогда  $\{i_1, \dots, i_k\}$  сопоставляется вектор  $\omega = (\omega_1, \dots, \omega_n)$ , все компоненты которого равны 0 кроме равных 1 компонент  $(\omega_{i_1}, \dots, \omega_{i_k})$ . Теоретические исследования свойств тупиковых тестов для случайных бинарных таблиц показали, что характеристические векторы для почти всех тупиковых тестов имеют асимптотически (при неограниченном возрастании размерности таблицы обучения) одну и ту же длину. Данный результат является обоснованием выбора в качестве системы опорных векторов всевозможных наборов, включающих фиксированное число признаков  $k$  или

$$\Omega_A = \{\omega : |\omega| = k\}$$

Оптимальное значение  $k$  находится в процессе обучения или задаётся экспертом. .

Другой часто используемой системе опорных множеств соответствует множество всех подмножеств  $\{1, \dots, n\}$  за исключением пустого множества. Иными словами в систему опорных множеств входит произвольный набор признаков или  $\Omega = \{\omega\} \setminus \omega_0$ , где  $\omega_0$  - вектор, все компоненты которого равны 0.

**Задание функции близости.** Пусть опорное множество  $\{i_1, \dots, i_k\}$  соответствует характеристическому вектору  $\omega$ . Фрагмент  $(x_{\mu i_1}, \dots, x_{\mu i_k})$  описания  $(x_{\mu 1}, \dots, x_{\mu n})$  объекта  $s_\mu$  называется  $\omega$ -частью объекта  $s_\mu$ .

Под функцией близости  $B_\omega(s_\mu, s_\nu)$  понимается функция от соответствующих  $\omega$ -частей сравниваемых объектов, принимающая значения 1 (объекты близки) или 0 (объекты удалены). Функции близости обычно задаются с помощью порговых параметров  $(\varepsilon_1, \dots, \varepsilon_n)$ , характеризующих близость объектов по отдельным признакам.

## Пример функций близости.

1)  $B_{\omega}(s_{\mu}, s_{\nu}) = 1$ , если при произвольном  $i \in \{1, \dots, n\}$ , при котором  $\omega_i = 1$ , всегда выполняется неравенство

$$|x_{\mu i} - x_{\nu i}| < \varepsilon_i.$$

$B_{\omega}(s_{\mu}, s_{\nu}) = 0$ , если существует такое  $i' \in \{1, \dots, n\}$ , что одновременно  $\omega_{i'} = 1$  и  $|x_{\mu i'} - x_{\nu i'}| > \varepsilon_{i'}$ .

2) Пусть  $\varepsilon$  - скалярный пороговый параметр. Функция  $B_{\omega}(s_{\mu}, s_{\nu}) = 1$ , если выполняется неравенство

$$\left[ \sum_{i=1}^n \omega_i |x_{\mu i} - x_{\nu i}| \right] < \varepsilon.$$

В противном случае  $B_{\omega}(s_{\mu}, s_{\nu}) = 0$ .

Важным элементом АВО является оценка близости распознаваемого объекта  $s_*$  к эталону  $s_\mu$  по заданной  $\omega$ - части. Данная оценка близости, которая будет обозначаться  $\Gamma_\omega(s_*, s_\mu)$ , формируется на основе введённых ранее функций близости и, возможно, дополнительных параметров.

$$A) \Gamma_\omega(s_*, s_\mu) = B_\omega(s_*, s_\mu).$$

$$B) \Gamma_\omega(s_*, s_\mu) = p_\omega B_\omega(s_*, s_\mu),$$

где  $p_\omega$ - параметр, характеризующий информативность опорного множества с характеристическим вектором  $\omega$ .

$$C) \Gamma_\omega(s_*, s_\mu) = \gamma_\mu \left( \sum_{j=1}^n p_j \omega_j \right) B_\omega(s_*, s_\mu),$$

где  $\gamma_\mu$  - параметр, характеризующий информативность эталонного объекта  $s_\mu$ , параметры  $p_1, \dots, p_n$  характеризуют информативность отдельных признаков.

Оценка объекта  $s_*$  за класс  $K_l$  при фиксированном  $\omega$ . Оценка объекта  $s_*$  за класс  $K_l$  при фиксированном характеристическом векторе  $\omega$  может вычисляться как среднее значение близости  $s_*$  к эталонным объектам из класса  $K_l$

$$\Gamma_{\omega}^l(s_*) = \frac{1}{m_l} \sum_{s_{\mu} \in K_l} \Gamma_{\omega}(s_*, s_{\mu}).$$

Общая оценка  $s_*$  за класс  $K_l$  вычисляется как сумма оценок  $\Gamma_{\omega}^l(s_*)$  по опорным множествам из системы  $\Omega_A$ :

$$\Gamma^l(s_*) = \sum_{\omega \in \Omega_A} \Gamma_{\omega}^l(s_*). \quad (1)$$

Наряду с формулой (1) используется формула

$$\Gamma^l(s_*) = \nu_l \sum_{\omega \in \Omega_A} \Gamma_{\omega}^l(s_*). \quad (2)$$

Использование взвешивающих параметров  $\nu_1, \dots, \nu_L$  позволяет регулировать доли правильно распознанных объектов  $K_1, \dots, K_L$ . Прямое вычисление оценок за классы по формулам (1) и (2) в случаях, когда в качестве систем опорных множеств используются наборы с фиксированным числом признаков или всевозможные наборы признаков, оказывается практически невозможным при сколь угодно высокой размерности признакового пространства из-за необходимости вычисления огромного числа значений функций близости. Однако при равенстве весов всех признаков существуют эффективные формулы для вычисления оценок по формуле (1). Предположим, что оценки близости распознаваемого объекта  $s_*$  к эталону  $s_\mu$  по заданной  $\omega$  - части вычисляются по формуле (А). Тогда оценка по формуле (1) принимает вид

$$\Gamma^l(s_*) = \sum_{\omega \in \Omega_A} \sum_{s_\mu \in K_l} B_{\omega(s_*, s_\mu)}$$

Рассмотрим сумму  $\sum_{\omega \in \Omega_A} B_{\omega}(s_*, s_{\mu})$ . Предположим, что общее число признаков, по которым объект  $s_*$  близок к объекту  $s_{\mu}$  равно  $d(s_*, s_{\mu})$ . Иными словами  $d(s_*, s_{\mu}) = |D(s_*, s_{\mu})|$ , где  $D(s_*, s_{\mu}) = \{i : |x_{*i} - x_{\mu i}| < \varepsilon_i\}$ . Очевидно функция близости  $B_{\omega}(s_*, s_{\mu}) = 1$  тогда и только тогда, когда опорное множество, задаваемое характеристическим вектором  $\omega$ , полностью входит в множество  $D(s_*, s_{\mu})$ . Во всех остальных случаях  $B_{\omega}(s_*, s_{\mu}) = 0$ .

Предположим, что система опорных множеств удовлетворяет условию  $\Omega_A = \{\omega : |\omega| = k\}$ . Очевидно, что число опорных множеств в  $\Omega_A$ , удовлетворяющих условию  $B_\omega(s_*, s_\mu) = 1$ , равно  $C_{d(s_*, s_\mu)}^k$ . Откуда следует, что  $\sum_{\omega \in \Omega_A} B_\omega(s_*, s_\mu) = C_{d(s_*, s_\mu)}^k$ . Следовательно оценка по формуле (1) может быть записана в виде

$$\Gamma^l(s_*) = \frac{1}{m_l} \sum_{s_\mu \in K_l} \gamma_\mu C_{d(s_*, s_\mu)}^k \quad (3)$$

Предположим, что система  $\Omega_A$  включает в себя всевозможные опорные множества. В этом случае число опорных множеств в  $\Omega_A$ , удовлетворяющих условию  $B_\omega(s_*, s_\mu) = 1$ , равно  $2^{d(s_*, s_\mu)} - 1$ . Следовательно оценка по формуле (1) может быть записана в виде

$$\Gamma^l(s_*) = \frac{1}{m_l} \sum_{s_\mu \in K_l} \gamma_\mu [2^{d(s_*, s_\mu)} - 1].$$

Для обучения алгоритмов АВО в общем случае может быть использован тот же самый подход, который используется для обучения в методе «Линейная машина». Предположим, что решается задача обучения алгоритмов для распознавания объектов, принадлежащих классам  $K_1, \dots, K_L$ . При правильного распознавания объекта  $s_i \in \tilde{S}_t \cap K_l$  должен выполняться блок неравенств

$$\Gamma^l(s_i) > \Gamma^{l'}(s_i), \quad (4)$$

где  $l' \in \{1, \dots, L\} \setminus \{l\}$ . При использовании для вычисления оценок формулы (3) блок (4) приводится к виду

$$\frac{1}{m_l} \sum_{s_\mu \in K_l} \gamma_\mu C_{d(s_*, s_\mu)}^k > \frac{1}{m_{l'}} \sum_{s_\nu \in K_{l'}} \gamma_\nu C_{d(s_*, s_\nu)}^k \quad (5)$$

где  $l' \in \{1, \dots, L\} \setminus \{l\}$ .

Обучение сводится к поиску такого набора весовых коэффициентов  $\gamma_1, \dots, \gamma_m$  при которых выполняется по возможности максимальное число блоков неравенств вида (4). Поиск максимальной оптимальных коэффициентов может производиться с использованием эвристического релаксационного метода, аналогичного тому, что был использован при обучении алгоритма «Линейная машина».