

## СПЕЦКУРС

### Логический анализ данных в распознавании (Logical data analysis in recognition)

*лектор д.ф.-м.н. Елена Всеволодовна Дюкова*

Спецкурс посвящён вопросам применения аппарата дискретной математики в задачах интеллектуального анализа данных. Излагаются общие принципы, лежащие в основе логического подхода к задачам машинного обучения. Описываются методы конструирования процедур классификации по прецедентам с использованием понятий теории булевых функций и теории покрытий булевых матриц. Рассматриваются основные модели логических процедур классификации, вопросы сложности их реализации и качества решения прикладных задач.

**Спецкурс для бакалавров 2-4 курсов ВМК МГУ им. М.В. Ломоносова.**

По спецкурсу издано учебное пособие:

<http://www.ccas.ru/frc/papers/djukova03mp.pdf>

## Лекция 6

## Использование аппарата логических функций для конструирования дискретных процедур распознавания в случае целочисленной информации

- Пусть  $E_k^n$ ,  $k \geq 2$ , - множество наборов вида  $(\alpha_1, \dots, \alpha_n)$ , где  $\alpha_i \in \{0, 1, \dots, k-1\}$ . Пусть переменная  $x$  принимает значения из множества  $\{0, 1, \dots, k-1\}$ ,  $\sigma \in \{0, 1, \dots, k-1\}$ . Положим

$$x^\sigma = \begin{cases} 1, & \text{если } x = \sigma, \\ 0, & \text{если } x \neq \sigma. \end{cases}$$

- *Элементарной конъюнкцией* (ЭК) над переменными  $x_1, \dots, x_n$  назовём функцию вида  $x_{j_1}^{\sigma_1} \& \dots \& x_{j_r}^{\sigma_r}$ , где  $\sigma_i \in \{0, 1, \dots, k-1\}$ ,  $x_{j_i} \in \{x_1, \dots, x_n\}$  при  $i = 1, 2, \dots, r$  и  $x_{j_q} \neq x_{j_t}$  при  $t, q \in \{1, 2, \dots, r\}$ ,  $t \neq q$ . Также, как и в случае  $k=2$ , для краткости знак  $\&$  опускается.

- Нетрудно видеть, что ЭК  $B = x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r}$  обращается в 1 на тех и только тех наборах  $(\alpha_1, \dots, \alpha_n)$ , для которых  $\alpha_{j_i} = \sigma_i$ ,  $i = 1, 2, \dots, r$ , т.е. ЭК принимает значение 1 на наборе из  $E_k^n$ , если каждый её сомножитель принимает значение 1 на этом наборе. Интервал истинности ЭК  $B$  будем обозначать через  $N_B$ .
- Пусть  $f(x_1, \dots, x_n)$  – функция, частично определённая на наборах из  $E_k^n$  и принимающая значения из  $\{0, 1\}$ .  $N_f$  и  $N_{\bar{f}}$  – соответственно множество единиц и множество нулей функции  $f$ .
- Определения почти допустимой, допустимой, неприводимой и максимальной конъюнкций, данные в лекции 4 для случая всюду определённой булевой функции, полностью переносятся на случай двузначной частичной функции  $f$ , определённой на наборах из  $E_k^n$ .

- Рассмотрим ситуацию, когда объекты из исследуемого множества  $M$  описаны признаками, каждый из которых принимает значения из множества  $\{0, 1, \dots, k - 1\}$ .
- Эл.кл.  $(\sigma, H)$ ,  $\sigma = (\sigma_1, \dots, \sigma_r)$ , порождённому набором признаков  $H = \{x_{j_1}, \dots, x_{j_r}\}$ , поставим в соответствие ЭК  $B = x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r}$ .
- Напомним (см. лекцию 3), что близость объекта  $S = (a_1, \dots, a_n)$  из  $M$  и элементарного классификатора  $(\sigma, H)$  оценивается величиной  $B(\sigma, S, H)$ , равной  $1$ , если  $a_{j_t} = \sigma_t$  при  $t = 1, 2, \dots, r$ , и равной  $0$  в противном случае.
- Очевидно,  $B(\sigma, S, H) = 1$  тогда и только тогда, когда  $(a_1, \dots, a_n) \in N_B$ .

- Рассмотрим основные модели логических классификаторов, а именно, алгоритмы голосования по представительным наборам, по антипредставительным наборам, по покрытиям классов и тестам. Покажем, что в каждой из этих моделей построение множества эл.кл. для класса  $K$  сводится к построению допустимых или максимальных конъюнкций двухзначной логической функции, частично или всюду определённой на наборах из  $E_k^n$ . Такая функция на обучающих объектах из  $K$  и  $\bar{K} = \{K_1, \dots, K_l\} \setminus K$  принимает разные значения и называется *характеристической функцией* класса  $K$ .
- Процедура голосования по построенной ЭК  $B$  заключается в проверке принадлежности набора  $(a_1, \dots, a_n)$  интервалу  $N_B$ .

## 1. Алгоритм голосования по представительным наборам

- В данном случае характеристическая функция класса  $K$  – частичная логическая функция  $f_K(x_1, \dots, x_n)$ , принимающая значение 1 на описаниях обучающих объектов из класса  $K$ , значение 0 на описаниях обучающих объектов из  $\bar{K}$  и не определённая на остальных наборах из  $E_K^n$ .
- Представительному набору класса  $K$  соответствует допустимая конъюнкция функции  $f_K$ , тупиковому представительному набору соответствует максимальная конъюнкция функции  $f_K$ .
- Допустимая (максимальная) конъюнкция  $B$  функции  $f_K$  голосует за принадлежность объекта  $S$  классу  $K$ , если  $(a_1, \dots, a_n) \in N_B$ .

## 2. Алгоритм голосования по антипредставительным наборам

- Характеристическая функция класса  $K$  – частичная логическая функция  $f_{\bar{K}}(x_1, \dots, x_n)$ , принимающая значение 0 на описаниях обучающих объектов из класса  $K$ , значение 1 на описаниях обучающих объектов из  $\bar{K}$  и не определённая на остальных наборах из  $E_k^n$ .
- Антипредставительному набору класса  $K$  соответствует допустимая конъюнкция функции  $f_{\bar{K}}$ , тупиковому антипредставительному набору соответствует максимальная конъюнкция функции  $f_{\bar{K}}$ .
- Допустимая (максимальная) конъюнкция  $B$  функции  $f_{\bar{K}}$  голосует за принадлежность объекта  $S$  классу  $K$ , если  $(a_1, \dots, a_n) \notin N_B$ .

### 3. Алгоритмы голосования по покрытиям класса

- Характеристическая функция класса  $K$  – всюду определённая логическая функция  $F_{\bar{K}}(x_1, \dots, x_n)$ , принимающая значение 0 на описаниях обучающих объектов из класса  $K$  и значение 1 на остальных наборах из  $E_k^n$ .
- Покрытие класса  $K$  соответствует допустимая конъюнкция функции  $F_{\bar{K}}$ , тупиковому покрытию – максимальная конъюнкция для  $F_{\bar{K}}$ .
- Допустимая (максимальная) конъюнкция  $B$  функции  $F_{\bar{K}}$  голосует за принадлежность объекта  $S$  классу  $K$ , если  $(a_1, \dots, a_n) \notin N_B$ .

#### 4. Алгоритм голосования по тестам

- Характеристическая функция класса  $K$  определяется так же, как и для алгоритмов голосования по представительным наборам.
- Множество всех почти допустимых конъюнкций функции  $f_K$ , порождаемое набором признаков  $H$ , обозначим через  $Q(H, K)$ .
- Очевидно, набор признаков  $H$  является тестом тогда и только тогда, когда для каждого  $t \in \{1, 2, \dots, l\}$ , каждая конъюнкция из  $Q(H, K_t)$  является допустимой для  $f_{K_t}$ .
- Очевидно также, что тест  $H = \{x_{j_1}, \dots, x_{j_r}\}$  является тупиковым тогда и только тогда, когда для каждого  $i$ ,  $i \in \{1, 2, \dots, r\}$ , в  $\{1, 2, \dots, l\}$  можно указать  $t_i$  такое, что  $Q(H \setminus \{x_i\}, K_{t_i})$  содержит конъюнкцию, не являющуюся допустимой для соответствующей характеристической функции.

- Построение требуемого множества конъюнкций (ДНФ, реализующей характеристическую функцию) может быть осуществлено на основе преобразования нормальных форм.
- Случай, когда  $k = 2$  рассмотрен в лекции 5. Описанные в лекции 5 алгоритмы построения сокращённой ДНФ булевой функции, заданной множеством нулей, могут быть обобщены на рассматриваемый нами общий случай ( $k \geq 2$ ).
- Пусть множество нулей характеристической функции  $f$  состоит из наборов  $(\beta_{11}, \dots, \beta_{1n}), (\beta_{21}, \dots, \beta_{2n}), \dots, (\beta_{u1}, \dots, \beta_{un})$ .
- Если  $k = 2$  и  $f$  всюду определена, то нужно построить КНФ  $K = D_1 \& \dots \& D_u$ , где  $D_i = x_1^{\overline{\beta_{i1}}} \vee \dots \vee x_n^{\overline{\beta_{in}}}$ ,  $i = 1, 2, \dots, u$ , реализующую функцию  $f$ . Затем преобразовать эту КНФ в (сокращённую) ДНФ функции  $f$ . Для этого можно воспользоваться одним из описанных в лекции 5 способов преобразования нормальных форм булевой функции.

- Построение требуемого множества конъюнкций (ДНФ, реализующей характеристическую функцию) может быть осуществлено на основе преобразования нормальных форм.
- Случай, когда  $k = 2$  рассмотрен в лекции 5. Описанные в лекции 5 алгоритмы построения сокращённой ДНФ булевой функции, заданной множеством нулей, могут быть обобщены на рассматриваемый нами общий случай ( $k \geq 2$ ).
- Пусть множество нулей характеристической функции  $f$  состоит из наборов  $(\beta_{11}, \dots, \beta_{1n}), (\beta_{21}, \dots, \beta_{2n}), \dots, (\beta_{u1}, \dots, \beta_{un})$ .
- Если  $k = 2$  и  $f$  всюду определена, то нужно построить КНФ  $K = D_1 \& \dots \& D_u$ , где  $D_i = x_1^{\overline{\beta_{i1}}} \vee \dots \vee x_n^{\overline{\beta_{in}}}$ ,  $i = 1, 2, \dots, u$ , реализующую функцию  $f$ . Затем преобразовать эту КНФ в (сокращённую) ДНФ функции  $f$ . Для этого можно воспользоваться одним из описанных в лекции 5 способов преобразования нормальных форм булевой функции.

- Если же  $k = 2$  и  $f$  не всюду определена, то КНФ  $K$  реализует всюду определённую булеву функцию  $F$  с тем же множеством нулей. В этом случае необходимо удалить из (сокращённой) ДНФ функции  $F$  те конъюнкции, которые не являются допустимыми для  $f$ .
- В случае  $k > 2$  вместо  $D_i = \overline{x_1^{\beta_{i1}}} \vee \dots \vee \overline{x_n^{\beta_{in}}}$  необходимо взять дизъюнкцию вида  $\overline{x_1^{\beta_{i1}}} \vee \dots \vee \overline{x_n^{\beta_{in}}}$  и воспользоваться равенством  $\overline{x^\beta} = \bigvee_{\alpha \neq \beta} x^\alpha$ . Тогда преобразуемая КНФ примет вид

$$D_1^* \& D_2^* \& \dots \& D_n^*,$$

где  $D_i^* = \bigvee_{\alpha \neq \beta_{i1}} x_1^\alpha \bigvee_{\alpha \neq \beta_{i2}} x_2^\alpha \bigvee \dots \bigvee_{\alpha \neq \beta_{in}} x_n^\alpha$ .

- После перемножения логических скобок в получившейся ДНФ следует сделать упрощения, пользуясь тем, что  $x^\alpha \cdot x^\beta = 0$  при  $\alpha \neq \beta$ ,  $x \cdot x = x$ ,  $x \vee x = x$ ,  $x \vee xx' = x$ . Остальные рассуждения полностью совпадают с рассуждениями, приведёнными в лекции 5.

- Построение требуемого множества тестов может быть осуществлено на основе преобразования КНФ, не содержащей отрицаний переменных (реализующей монотонную булеву функцию), в ДНФ.
- Действительно, рассмотрим множество  $P$  всех неупорядоченных пар наборов  $(\alpha, \beta)$  таких, что  $\alpha \in N_{f_{k_u}}$ ,  $\beta \in N_{f_{k_v}}$ ,  $u \neq v$ ,  $u, v \in \{1, 2, \dots, l\}$ . Для каждой пары  $(\alpha, \beta)$  из  $P$  построим дизъюнкцию  $D(\alpha, \beta) = x_{p_1} \vee \dots \vee x_{p_q}$ , где  $p_1, \dots, p_q$  - номера разрядов, в которых различаются наборы  $\alpha$  и  $\beta$ .
- Нетрудно убедиться в том, что для построения множества тестов нужно преобразовать КНФ  $\&_{(\alpha, \beta) \in P} D(\alpha, \beta)$  в ДНФ, состоящую из допустимых конъюнкций функции, реализуемой этой КНФ. Для построения множества тупиковых тестов нужно преобразовать КНФ  $\&_{(\alpha, \beta) \in P} D(\alpha, \beta)$  в ДНФ, состоящую из максимальных конъюнкций функции, реализуемой этой КНФ.

## УПРАЖНЕНИЯ

1. Пусть двузначная логическая функция  $F(x_1, \dots, x_n)$  определена на  $E_K^n$  и принимает значение 0 на описаниях обучающих объектов из класса  $K$ , а значение 1 на остальных наборах из  $E_K^n$ . Пусть конъюнкция  $x_{j_1}^{\sigma_1} \& \dots \& x_{j_r}^{\sigma_r}$  является неприводимой для  $F$ . Является ли эл.кл.  $(\sigma, H)$ , где  $\sigma = (\sigma_1, \dots, \sigma_r)$ ,  $H = (x_{j_1}, \dots, x_{j_r})$ , тупиковым покрытием класса  $K$  ?