

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 7

Поиск элементарных классификаторов на основе построения покрытий булевой и целочисленной матриц. Связь задач построения покрытий булевой матрицы и преобразования нормальных форм булевой функции

- При реализации логических классификаторов для поиска элементарных классификаторов чаще используются построения, в основе которых лежит поиск покрытий булевой или целочисленной матрицы.
- Пусть L – булева матрица. Набор столбцов H матрицы L называется *покрытием*, если каждая строка матрицы L в пересечении хотя бы с одним из столбцов, входящих в H , дает **1**. Покрытие называется *неприводимым*, если никакое его собственное подмножество не является покрытием.
- При построении неприводимых покрытий обычно используется следующий *критерий*.

- **Утверждение 1.** Набор H из r различных столбцов матрицы L является неприводимым покрытием, если выполнены два условия: 1) подматрица L^H матрицы L образованная столбцами набора H , не содержит строк вида $(0, 0, \dots, 0)$; 2) L^H содержит каждую из строк вида каждую из строк вида $(1, 0, \dots, 0), (0, 1, \dots, 0), \dots, (0, 0, \dots, 1)$, т.е. L^H содержит единичную подматрицу порядка r .
- Первое условие называется условием *покрываемости*, а второе условием *совместимости*.
- Как возникают задачи построения покрытий и неприводимых покрытий при построении (тупиковых) тестов и (тупиковых) представительных наборов?

- Материал обучения обычно представляется в виде таблицы T с n столбцами (таблицы обучения), в которой столбец с номером j соответствует признаку x_j , а каждая строка есть набор значений признаков, описывающих один из обучающих объектов. Для нахождения искомого множества элементарных классификаторов строится специальная булева матрица (матрица сравнения таблицы T). Обозначим её L_T . Каждая строка этой матрицы образуется в результате сравнения пары строк таблицы T , описывающих объекты из разных классов. При этом в столбце матрицы L_T с номером j ставится 1 , если сравниваемые строки различаются в разряде с номером j , и 0 в противном случае.
- Обозначим через $L_T^{(i)}$, $i = 1, \dots, m$, подматрицу матрицы L_T , которая образована сравнением обучающего объекта S_i со всеми обучающими объектами, не принадлежащими тому же классу, что и объект S_i .

- Очевидными являются следующие два утверждения.
- **Утверждение 2.** Набор признаков $\{x_{j_1}, \dots, x_{j_r}\}$ является (тупиковым) тестом тогда и только тогда, когда набор столбцов матрицы L_T с номерами j_1, \dots, j_r является (неприводимым) покрытием.
- **Утверждение 3.** Эл.кл. $(\sigma, \{x_{j_1}, \dots, x_{j_r}\})$, является (тупиковым) представительным набором для класса K тогда и только тогда, когда $B(\sigma, S_i, \{x_{j_1}, \dots, x_{j_r}\}) = \mathbf{1}$ для некоторого обучающего объекта S_i из K и набор столбцов матрицы $L_T^{(i)}$ с номерами j_1, \dots, j_r является (неприводимым) покрытием.
- Задача построения неприводимых покрытий булевой матрицы L размера $u \times n$ может быть сформулирована как задача преобразования КНФ монотонной булевой функции в сокращенную ДНФ.

- Действительно, строке с номером i поставим в соответствие дизъюнкцию $D_i = x_{p_1} \vee \dots \vee x_{p_q}$, где p_1, \dots, p_q – номера тех столбцов, которые в пересечении с этой строкой дают **1**. Пусть f -монотонная булева функция реализуемая КНФ $D_1 \& \dots \& D_u$.
- Пользуясь утверждением 2 из лекции 5 нетрудно доказать
- **Утверждение 4.** ЭК $x_{j_1} \dots x_{j_r}$ является допустимой для f тогда и только тогда, когда набор столбцов H матрицы L с номерами j_1, \dots, j_r является покрытием.
- Пользуясь утверждением 4 из лекции 5 нетрудно доказать
- **Утверждение 5.** ЭК $x_{j_1} \dots x_{j_r}$ является неприводимой для f тогда и только тогда, когда набор столбцов матрицы L с номерами j_1, \dots, j_r содержит единичную подматрицу порядка r .

- Из утверждений 4 и 5 следует
- **Утверждение 6.** ЭК $x_{j_1} \dots x_{j_r}$ является максимальной для f тогда и только тогда, когда набор столбцов матрицы L с номерами j_1, \dots, j_r является неприводимым покрытием.
- Из последних трех утверждений следует, что алгоритмы построения неприводимых покрытий булевой матрицы можно применять для построения максимальных конъюнкций монотонной булевой функции, заданной КНФ, и наоборот. В теории алгоритмической сложности рассматриваемая задача преобразования нормальных форм монотонной булевой функции называется *монотонной дуализацией* и относится к числу труднорешаемых дискретных задач. О существующих подходах к её решению будет рассказано в следующей лекции. В силу вычислительной сложности монотонной дуализации на практике отказываются от поиска всех тупиковых тестов и тупиковых представительных наборов. В основном используются модели с ограничением на длину эл.кл. и стохастические алгоритмы, в которых строится достаточно представительная выборка из искомого множества элементарных классификаторов.

- Можно избежать построения вспомогательной булевой матрицы L_T , если ввести понятие покрытия более общего вида.
- Действительно, пусть L – целочисленная матрица размера $u \times n$ с элементами из $\{0, 1, \dots, k - 1\}$, $k \geq 2$; E_k^r , $r \leq n$, – множество наборов вида $\sigma = (\sigma_1, \dots, \sigma_r)$, где $\sigma_i \in \{0, 1, \dots, k - 1\}$.
- Пусть далее $\sigma \in E_k^r$, $\sigma = (\sigma_1, \dots, \sigma_r)$.
- Набор H из r различных столбцов матрицы L называется σ -покрытием, если в подматрице L^H матрицы L , образованной столбцами набора H , нет строки $(\sigma_1, \dots, \sigma_r)$. Набор H из r различных столбцов матрицы L , являющийся σ -покрытием, называется *тупиковым σ -покрытием*, если L^H содержит подматрицу, имеющую с точностью до перестановки строк вид

$$\begin{pmatrix} \beta_1 & \sigma_2 & \sigma_3 & \dots & \sigma_{r-1} & \sigma_r \\ \sigma_1 & \beta_2 & \sigma_3 & \dots & \sigma_{r-1} & \sigma_r \\ \sigma_1 & \sigma_2 & \sigma_3 & \dots & \sigma_{r-1} & \beta_r \end{pmatrix},$$

где $\beta_p \neq \sigma_p$ при $p = 1, 2, \dots, r$. Такая подматрица называется σ -подматрицей.

- Если $k = 2$ и $\sigma = (0, \dots, 0)$, то понятие (тупикового) σ -покрытия совпадает с понятием (неприводимого) покрытия. Аналогом единичной подматрицы в этом случае является σ -подматрица.
- Таблицу обучения T можно рассматривать как пару матриц L_1 и L_2 , где L_1 - матрица, состоящая из описаний обучающих объектов из класса K , L_2 - матрица, состоящая из описаний остальных обучающих объектов. Тогда, очевидно, эл.кл. $(\sigma, \{x_{j_1}, \dots, x_{j_r}\})$, $\sigma = (\sigma_1, \dots, \sigma_r)$, является (тупиковым) представительным набором для класса K , если набор столбцов матрицы L_1 с номерами j_1, \dots, j_r не является $(\sigma_1, \dots, \sigma_r)$ -покрытием, а набор столбцов матрицы L_2 с номерами j_1, \dots, j_r является (тупиковым) $(\sigma_1, \dots, \sigma_r)$ -покрытием.
- Нетрудно видеть, что эл.кл. $(\sigma, \{x_{j_1}, \dots, x_{j_r}\})$, $\sigma = (\sigma_1, \dots, \sigma_r)$, является (тупиковым) покрытием класса K , если набор столбцов матрицы L_1 с номерами j_1, \dots, j_r является (тупиковым) $(\sigma_1, \dots, \sigma_r)$ -покрытием.

• Нетрудно также видеть, что эл.кл. вида $(\sigma, \{x_{j_1}, \dots, x_{j_r}\})$, $\sigma = (\sigma_1, \dots, \sigma_r)$, является (тупиковым) антипредставительным набором для класса K , если набор столбцов матрицы L_2 с номерами j_1, \dots, j_r не является $(\sigma_1, \dots, \sigma_r)$ -покрытием, а набор столбцов матрицы L_1 с номерами j_1, \dots, j_r является (тупиковым) $(\sigma_1, \dots, \sigma_r)$ -покрытием.

• Через $R(\sigma)$ обозначим множество наборов $(\beta_1, \dots, \beta_r)$ в E_k^r таких, что $\beta_j \neq \sigma_j$ при $j = 1, 2, \dots, r$. Набор столбцов H матрицы L назовем $R(\sigma)$ -покрытием, если в подматрице L^H матрицы L , образованной столбцами набора H , нет ни одной строки из $R(\sigma)$. Набор столбцов H матрицы L , являющийся $R(\sigma)$ -покрытием, назовем *тупиковым $R(\sigma)$ -покрытием*, если L^H содержит подматрицу, имеющую с точностью до перестановки строк вид

$$\begin{pmatrix} \sigma_1 & \beta_2 & \beta_3 & \dots & \beta_{r-1} & \beta_r \\ \beta_1 & \sigma_2 & \beta_3 & \dots & \beta_{r-1} & \beta_r \\ & & & \dots & & \\ \beta_1 & \beta_2 & \beta_3 & \dots & \beta_{r-1} & \sigma_r \end{pmatrix},$$

где $\beta_p \neq \sigma_p$ при $p = 1, 2, \dots, r$. Такая подматрица называется $R(\sigma)$ -подматрицей.

- Если $k = 2$ и $\sigma = (1, \dots, 1)$, то понятие (тупикового) $R(\sigma)$ -покрытия совпадает с понятием (неприводимого) покрытия. Аналогом единичной подматрицы является $R(\sigma)$ -подматрица.
- Пусть далее L – булева матрица размера $u \times n$, $\sigma \in E_2^r$.
- Через $C(L, \sigma)$ и $B(L, \sigma)$ обозначим соответственно совокупность всех σ -покрытий и всех тупиковых σ -покрытий матрицы L . Положим
- $C(L) = \bigcup_{r=1}^n \bigcup_{\sigma \in E_2^r} C(L, \sigma)$, $B(L) = \bigcup_{r=1}^n \bigcup_{\sigma \in E_2^r} B(L, \sigma)$.
- Связь между задачами построения множеств $C(L)$, $B(L)$ матрицы L и задачей преобразования нормальных форм булевой функции устанавливается следующим образом.
- Пусть $(\sigma_{i1}, \dots, \sigma_{in})$ – строка матрицы L с номером i , $i \in \{1, 2, \dots, u\}$. Этой строке ставится в соответствие дизъюнкция $D_i^{(1)} = x_1^{\overline{\sigma_{i1}}} \vee \dots \vee x_n^{\overline{\sigma_{in}}}$

- Пусть F_1 – булева функция, реализуемая КНФ $D_1^{(1)} \& \dots \& D_u^{(1)}$. Используя утверждения 2 и 4 из лекции 5 нетрудно доказать приведенные ниже утверждения 7, 8, из которых сразу следует утверждение 9.
- **Утверждение 7.** ЭК $x_{j_1}^{\sigma_1} \& \dots \& x_{j_r}^{\sigma_r}$ является допустимой для F_1 тогда и только тогда, когда набор столбцов матрицы L с номерами j_1, \dots, j_r является $(\sigma_1, \dots, \sigma_r)$ -покрытием.
- **Утверждение 8.** ЭК $x_{j_1}^{\sigma_1} \& \dots \& x_{j_r}^{\sigma_r}$ является неприводимой для F_1 тогда и только тогда, когда набор столбцов матрицы L с номерами j_1, \dots, j_r содержит $(\sigma_1, \dots, \sigma_r)$ -подматрицу.
- **Утверждение 9.** ЭК $x_{j_1}^{\sigma_1} \& \dots \& x_{j_r}^{\sigma_r}$ является максимальной для F_1 тогда и только тогда, когда набор столбцов матрицы L с номерами j_1, \dots, j_r является тупиковым $(\sigma_1, \dots, \sigma_r)$ -покрытием.

- Если строке матрицы L с номером i , $i \in \{1, 2, \dots, u\}$, поставить в соответствие дизъюнкцию $D_i^{(2)} = x_1^{\sigma_{i1}} \vee \dots \vee x_n^{\sigma_{in}}$ и рассмотреть булеву функцию F_2 , реализуемую КНФ $D_1^{(2)} \& \dots \& D_u^{(2)}$, то задачи построения допустимых и максимальных конъюнкций функции F_2 могут быть сформулированы соответственно как задачи построения $R(\sigma)$ -покрытий и тупиковых $R(\sigma)$ -покрытий матрицы L .
- Таким образом, алгоритмы построения покрытий общего вида булевой матрицы можно применять для построения ДНФ булевой функции, заданной КНФ, и наоборот. Аналогичное утверждение имеет место и для целочисленной матрицы. Алгоритмы построения покрытий общего вида целочисленной матрицы можно применять для ДНФ двузначной логической функции k -значной логики, заданной КНФ.

УПРАЖНЕНИЯ

- Пусть две строки булевой матрицы L имеют вид (b_1, \dots, b_n) и (c_1, \dots, c_n) , где $c_j \geq b_j$ при $j = 1, 2, \dots, n$. Будем говорить, что вторая строка охватывает первую. Показать, что при удалении охватывающих строк из матрицы L множество её покрытий не меняется.
- 2. Пусть булева F функция задана КНФ $(x_1 \vee x_2)(x_2 \vee x_3)(x_1 \vee x_3)$. Построить сокращенную ДНФ функции F сведением задачи к задаче построения неприводимых покрытий булевой матрицы.
- 3. Пусть булева F функция задана множеством наборов, на которых она равна 0 , а именно, $N_{\bar{F}} = \{(1, 1, 1), (0, 1, 1), (1, 0, 1)\}$. Построить сокращенную ДНФ функции F двумя способами: а) сведением к задаче построения тупиковых σ -покрытий булевой матрицы; б) сведением к задаче построения тупиковых $R(\sigma)$ -покрытий булевой матрицы.
- 4. Доказать утверждения 4 – 9.