

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 10

Алгоритмы классификации на основе решающих деревьев

- Решающие деревья – это один из наиболее популярных инструментов для решения задач классификации по прецедентам.
- Для простоты рассмотрим задачу распознавания с бинарными признаками x_1, \dots, x_n , с непересекающимися классами K_1, \dots, K_l и обучающими объектами S_1, \dots, S_m .
- Ориентированное бинарное дерево называется бинарным решающим деревом (БРД), если:
 - Каждой внутренней вершине дерева сопоставлен один из признаков x_1, \dots, x_n
 - Каждой висячей вершине сопоставлен один из классов K_1, \dots, K_l
 - Каждая дуга помечена одним из чисел 0 или 1, причем из каждой внутренней вершины выходят две дуги, помеченные разными числами.

- Нетрудно видеть, что ветви БРД с вершинами x_{j_1}, \dots, x_{j_r} можно сопоставить ЭК над переменными x_1, \dots, x_n вида $x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r}$, где σ_i – метка дуги, выходящей из вершины x_{j_i} , $i = 1, 2, \dots, r$. С другой стороны, каждой висячей вершине, а следовательно, каждой ветви БРД сопоставлен один из классов K_1, \dots, K_l . Таким образом, каждой ветви БРД соответствует пара (B, K) , где B – ЭК над переменными x_1, \dots, x_n , $K \in \{K_1, \dots, K_l\}$.
- Пусть БРД имеет μ ветвей и этим ветвям соответствуют пары $(B_1, K_{j_1}), \dots, (B_\mu, K_{j_\mu})$, N_{B_i} – интервал истинности ЭК B_i , $i \in \{1, 2, \dots, \mu\}$, S – распознаваемый объект. Если среди конъюнкций B_1, \dots, B_μ существует конъюнкция B_i такая, что $S \in N_{B_i}$, то объект S относится к классу K_{j_i} . В противном случае происходит отказ от распознавания объекта S .

- Пусть некоторой ветви дерева сопоставлена пара (B, K) , где $B = x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r}$. Рассматриваемая ветвь называется корректной, если наборы из $\{S_1, \dots, S_m\} \cap N_B$ являются описаниями объектов класса K . В этом случае интервал истинности N_B конъюнкции B называется допустимым. Бинарное решающее дерево, у которого каждая ветвь корректная, называется корректным. Нетрудно видеть, что в корректном дереве ветвь (B, K) с конъюнкцией $B = x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r}$ порождает представительный набор класса K вида (σ, H) , $\sigma = (\sigma_1, \dots, \sigma_r)$, $H = (x_{j_1}, \dots, x_{j_r})$.
- Синтез решающего дерева представляет собой итерационный процесс. Как правило, для построения очередной вершины дерева выбирается признак, наилучшим образом удовлетворяющий некоторому критерию ветвления. По значениям этого признака и осуществляется ветвление, далее указанная процедура повторяется для каждой из ветвей. Описанный подход обуславливает основные достоинства метода, а именно, решающее правило строится быстро, получается достаточно простым и хорошо интерпретируемым.

- Одним из наиболее известных алгоритмов классификации на основе БРД является **алгоритм построения допустимого разбиения (далее АDR)**.
- Критерий ветвления в АDR представляет собой набор условий с разным приоритетом. При выборе очередной вершины сначала проверяется условие с наибольшим приоритетом. Ищется признак с наименьшим номером, для которого это условие выполняется. Если ни один признак не удовлетворяет рассматриваемому условию, то проверяется следующее по порядку условие с более низким приоритетом. Однако, в случае, когда два или более признака удовлетворяют рассматриваемому критерию ветвления в равной или почти равной мере, выбор одного из этих признаков осуществляется практически случайным образом. При этом в зависимости от выбранного признака построенные деревья могут существенно отличаться как по составу используемых признаков, так и по своим распознающим качествам.

- Алгоритм АДР строит разбиение множества E^n всех двоичных наборов длины n на непересекающиеся интервалы, такие, что в каждом из них находятся описания обучающих объектов только из одного класса.
- На первом шаге множество E^n разбивается на два непересекающихся интервала N_1^1 и N_2^2 ранга $n - 1$ таких, что $N_1^1 \cup N_2^2 = E^n$ и можно указать пару наборов из E^n , принадлежащих разным интервалам и описывающих обучающие объекты из разных классов. Интервалы N_1^1 и N_2^2 являются интервалами истинности конъюнкций x_{j_1} и \bar{x}_{j_1} ранга 1.
- Пусть на шаге i построено разбиение E^n на некоторое множество интервалов. Если в этом разбиении не существует интервала, содержащего наборы, описывающие объекты из разных классов, то алгоритм заканчивает свою работу. В противном случае в построенном разбиении выбирается интервал, содержащий наборы, описывающие объекты из разных классов. Пусть $n - r$ – ранг выбранного интервала и пусть ему соответствует конъюнкция $x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r}$. Тогда этот интервал разбивается на два непересекающихся интервала N_1^{i+1} и N_2^{i+1} ранга $n - r - 1$

• Пусть $n - r$ – ранг выбранного интервала и пусть ему соответствует конъюнкция $x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r}$. Тогда этот интервал разбивается на два непересекающихся интервала N_1^{i+1} и N_2^{i+1} ранга $n - r - 1$, таких, что можно указать пару наборов $\tilde{\alpha}$ и $\tilde{\beta}$, $\tilde{\alpha} \in N_1^{i+1}$, $\tilde{\beta} \in N_2^{i+1}$, которые описывают объекты из разных классов. Интервалами N_1^{i+1} и N_2^{i+1} соответствуют конъюнкции вида $x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r} \bar{x}_{j_{r+1}}$ и $x_{j_1}^{\sigma_1} \dots x_{j_r}^{\sigma_r} x_{j_{r+1}}$.

• Нетрудно видеть, что в процессе работы описанного выше алгоритма строится множество допустимых интервалов N_1, \dots, N_μ таких, что $\bigcup_{q=1}^{\mu} N_q = E^n$ и $\bigcap_{q=1}^{\mu} N_q = \emptyset$. Множество интервалов N_1, \dots, N_μ называется *допустимым разбиением*.

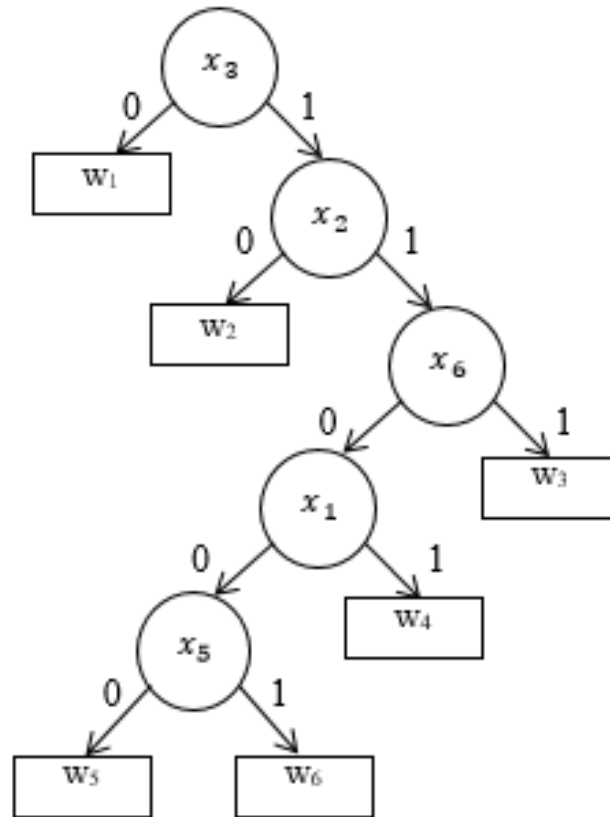
- Статистические оценки надежности распознавания БРД на контрольной выборке говорят о том, что наиболее надежным является решающее дерево с наименьшим числом листьев, а среди БРД, имеющих наименьшее число листьев, наиболее надежными являются те, у которых объекты обучающей выборки наиболее равномерно распределены по интервалам соответствующего разбиения. Известно, что задача построения корректного БРД с минимальным числом листьев является NP-полной. Поэтому при выборе признака для ветвления используются эвристические соображения, называемые критерием выбора признака. Например, при выборе признака последовательно проверяются условия полной, частичной и максимальной отделимости, приведенные ниже.
- Пусть для разбиения интервала N на пару интервалов N_0 и N_1 выбирается признак x_j . Имеем $N = N_0 \cup N_1$, $N_0 \cap N_1 = \emptyset$. Обозначим через $A_0(j)$ и $A_1(j)$ соответственно множество $\{S_1, \dots, S_m\} \cap N_0$ и $\{S_1, \dots, S_m\} \cap N_1$ (здесь и далее $\{S_1, \dots, S_m\}$ - обучающая выборка).

- Будем говорить, что признак x_j удовлетворяет *условию полной отделимости*, если каждое из множеств $A_0(j)$ и $A_1(j)$ содержит описания объектов только одного класса. Будем говорить, что признак x_j удовлетворяет *условию частичной отделимости*, если одно из множеств $A_0(j)$ или $A_1(j)$ содержит описания объектов только одного класса, а другое этим свойством не обладает. Пусть $D(u)$, $u \in 1, 2, \dots, n$, – число пар объектов вида (S_{i_1}, S_{i_2}) таких, что S_{i_1}, S_{i_2} – обучающие объекты из разных классов и $S_{i_1} \in A_0(u), S_{i_2} \in A_1(u)$. Будем говорить, что признак x_j удовлетворяет *условию максимальной отделимости*, если $D(j) = \max_{u \in \{1, \dots, n\}} D(u)$.
- Рассмотрим вариант алгоритма построения допустимого разбиения с использованием перечисленных выше условий. Данный алгоритм является рекурсивным. Положим на первом шаге $T^* = T_{mn}$, где T_{mn} – таблица из m строк и n столбцов, в которой строка с номером i , $i \in \{1, \dots, m\}$, представляет собой описание обучающего объекта S_i (таблица обучения). Далее на каждом шаге рекурсии выполняется следующая последовательность действий.

- Если все строки в T^* описывают объекты одного класса, то создается висячая вершина с меткой этого класса и рекурсия останавливается. В противном случае вызывается процедура выбора признака по таблице T^* .
- Пусть в результате выбран признак x_j . Создается внутренняя вершина с меткой x_j и для каждого значения $a \in \{0, 1\}$ признака x_j делается следующее. Создается дуга с меткой a , исходящая из построенной вершины. Строится подматрица $T^*(j, a)$ матрицы T^* путём удаления столбца, соответствующего признаку x_j , и удаления строк, которые в пересечении с этим столбцом дают a . Далее полагается $T^* = T^*(j, a)$ и производится следующий шаг рекурсии.

- Процедура выбора признака по таблице T^* представляет собой следующую последовательность шагов.
- Шаг 1. Для каждого признака проверяется условие полной отделимости. Если находятся признаки, для которых оно выполняется, то среди них выбирается первый по порядку признак и процедура останавливается.
- Шаг 2. Для каждого признака проверяется условие частичной отделимости. Если находятся признаки, для которых оно выполняется, то среди них выбирается первый по порядку признак и процедура останавливается.
- Шаг 3. Выбирается первый по порядку признак, для которого выполняется условие максимальной отделимости.

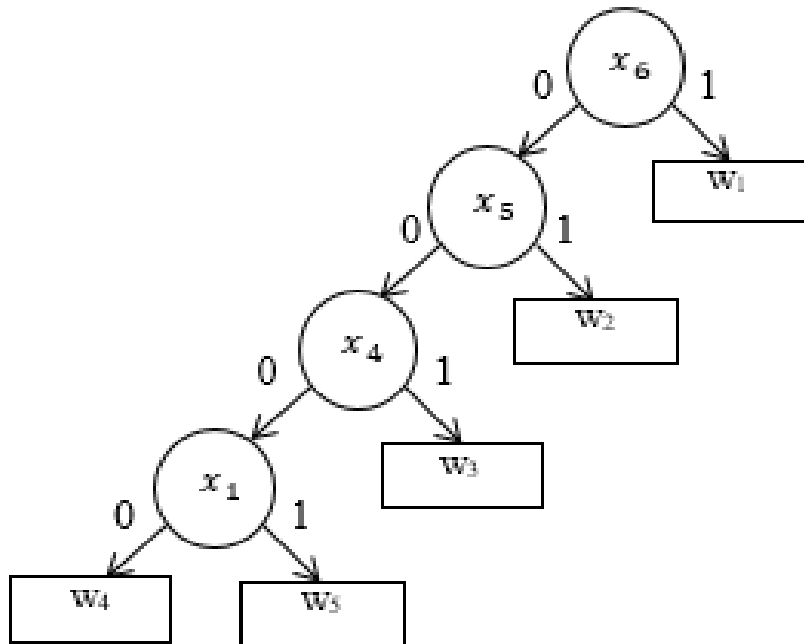
- Рассмотрим пример. Пусть обучающая выборка состоит из шести объектов и имеет вид:
 - $(1, 0, 0, 0, 1, 0)$,
 - $(0, 1, 1, 1, 1, 0)$,
 - $(1, 1, 1, 0, 0, 0)$,
 - $(1, 0, 1, 1, 0, 0)$,
 - $(1, 1, 1, 0, 1, 1)$,
 - $(0, 1, 1, 0, 0, 0)$.
- Пусть первые три объекта принадлежат классу первому классу, а последующие три объекта принадлежат второму классу.
- Ниже на рисунках 1 и 2 приведены различные БРД, построенные с помощью алгоритма АДР.



$$w_1 = (\bar{x}_3, K_1), w_2 = (\bar{x}_2 x_3, K_2), w_3 = (x_2 x_3 x_6, K_2), w_4 = (x_1 x_2 x_3 x_6, K_1),$$

$$w_5 = (\bar{x}_1 x_2 x_3 \bar{x}_5 x_6, K_2), w_6 = (\bar{x}_1 x_2 x_3 x_5 x_6, K_1)$$

Рис. 1



$$w_1 = (x_6, K_2), w_2 = (x_5 \bar{x}_6, K_1), w_3 = (x_4 \bar{x}_5 \bar{x}_6, K_2), w_4 = (\bar{x}_1 \bar{x}_4 \bar{x}_5 \bar{x}_6, K_2),$$

$$w_5 = (x_1 \bar{x}_4 \bar{x}_5 \bar{x}_6, K_1)$$

Рис. 2

- Как видно из приведенного выше описания алгоритма АДР, если проверяемому условию удовлетворяет более одного признака, то выбор одного из этих признаков происходит фактически случайно. Чтобы исправить указанный недостаток предлагается несколько модифицировать решающее дерево. Для этого введем второй тип внутренних вершин, которым соответствует набор равноценных признаков, т.е. признаков, удовлетворяющих критерию ветвления в равной или почти равной мере.
- Пусть внутренней вершине второго типа v , называемой далее полной вершиной, соответствует набор признаков $\{x_{j_1}, \dots, x_{j_q}\}$, $1 \leq q \leq n$, тогда из этой вершины исходит q дуг t_1, \dots, t_q , каждая из которых связывает вершину v с вершиной первого типа. Причем, дуга t_u , $u = 1, 2, \dots, q$, входит в вершину первого типа с меткой x_{j_u} . Полученную конструкцию будем называть полным решающим деревом.

- Заметим, что в полном решающем дереве, также как и в обычном, каждой висячей вершине соответствует некоторая конъюнкция. Однако в полном решающем дереве описание распознаваемого объекта может попасть в интервалы истинности конъюнкций, соответствующих разным ветвям дерева. В случае, если таким ветвям дерева соответствуют разные классы, выбор класса осуществляется простым голосование, т.е. объект зачисляется в тот класс, который соответствует большинству из указанных ветвей.
- Алгоритм построения полного решающего дерева (ПРД) также является рекурсивным. Обозначим через T^* матрицу, рассматриваемую на текущем шаге алгоритма, на первом шаге $T^* = T_{mn}$. Шаг рекурсии представляет собой следующую последовательность действий.

- Если все строки в таблице T^* описывают объекты одного класса, то создается висячая вершина с меткой этого класса и рекурсия останавливается. В противном случае вызывается процедура выбора в таблице T^* набора равноценных признаков для ветвления. Пусть в результате выбран набор признаков $\{x_{j_1}, \dots, x_{j_q}\}$. Тогда создается полная вершина с меткой $\{x_{j_1}, \dots, x_{j_q}\}$ и для каждого признака $x_{j_u}, u \in \{1, 2, \dots, q\}$, строится вершина с меткой x_{j_u} . Для каждого значения $a \in \{0, 1\}$ признака x_{j_u} создается дуга, исходящая из вершины x_{j_u} . Строится подматрица $T^*(j_u, a)$, которая образуется удалением столбца матрицы T^* , соответствующего признаку с номером j_u , и строк, в которых $a_{ij_u} \neq a, i = 1, 2, \dots, m$. Далее для рассматриваемой ветви полагается $T^* = T^*(j_u, a)$ и производится следующий шаг рекурсии.
- Процедура выбора по таблице T^* набора признаков для ветвления работает аналогично процедуре выбора признака для ветвления в классическом БДР.

- Описанный алгоритм также как АДР легко обобщается на случай, когда признаки целочисленные. Сравнение этого алгоритма с АДР, а также с широко используемым алгоритмом построения решающего дерева, известным под названием C4.5, показало целесообразность использования идеи полного решающего дерева. Результаты тестирования алгоритмов описаны в работе:
- *Djukova E.V., Peskov N.V. A Classification Algorithm Based on the Complete Decision Tree // J. Pattern Recognition and Image Analysis, 2007. Vol. 17. No. 3, pp. 363–367.*

- Затем была поставлена задача модифицировать на основе ПРД алгоритм С4.5 (Генрихов И.Е., Дюкова Е.В. *Классификация на основе полных решающих деревьев. Ж. вычисл. матем. и матем. физ.* 2012. Т. 52, № 4. С. 750-761). Впервые были построены и исследованы классификаторы на основе ПРД с энтропийным критерием ветвления, использующие взвешенное голосование по голосующим ветвям дерева (вместо голосования по большинству). С их помощью были эффективно решены вопросы обработки вещественнозначной информации. Построенные классификаторы существенно отличаются от алгоритма С4.5.
- Теоретическое и экспериментальное исследование обобщающей способности классификаторов на основе ПРД, т.е. способности этих классификаторов правильно классифицировать объекты, не вошедшие в обучающую выборку, свидетельствует о том, что обобщающая способность ПРД выше обобщающей способности классического РД и не уступает другим современным конструкциям РД, среди которых Random Forest.

УПРАЖНЕНИЯ

- 1. Классифицировать объекты $(1, 1, 1, 1, 1, 1)$ и $(0, 0, 0, 0, 0, 0)$, используя решающие деревья на рисунках 1 и 2.
- 2. Для рассмотренного в лекции примера построить полное решающее дерево.