

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 2

Принципы конструирования классических логических процедур распознавания и основные модели

- При решении прикладных задач достоверная информация о структуре множества M , как правило, отсутствует, поэтому при построении алгоритма распознавания мы не можем гарантировать качество работы этого алгоритма на новых объектах, отличных от прецедентов. Однако, если обучающие примеры достаточно характерны для исследуемого множества объектов, то алгоритм, редко ошибающийся на обучении, будет давать неплохие результаты и на неизвестных (не входящих в обучающую выборку) объектах. В связи с этим большое внимание уделяется проблеме корректности распознающих алгоритмов. Алгоритм называется *корректным*, если все объекты из обучающей выборки он распознает правильно.
- Будем предполагать, что описания объектов даны в виде значений целочисленных признаков и множество допустимых значений каждого признака ограничено. Рассмотрим простейший пример корректного распознающего алгоритма.

- Распознаваемый объект S будем сравнивать с каждым из объектов обучения S_1, \dots, S_m . В случае, если описание объекта S совпадает с описанием обучающего объекта S_i , то объект S будем относить к тому классу, которому принадлежит объект S_i , в противном случае будем считать, что алгоритм отказывается от распознавания. Нетрудно видеть, что описанный алгоритм является корректным, однако он не сможет распознать ни один объект, описание которого не совпадает с описанием ни одного из обучающих объектов.
- Очевидно, что требование полного совпадения описаний распознаваемого объекта и одного из обучающих объектов является слишком осторожным. Анализ прикладных задач свидетельствует о том, что близость объектов и их принадлежность одному классу можно оценивать на основе сравнения отдельных фрагментов их описаний. Поэтому возникает вопрос, как выбирать такие фрагменты, по которым будут сравниваться объекты. Один из вариантов ответа на данный вопрос используется в модели *алгоритмов вычисления оценок (АВО)*.

- Введем следующие обозначения. Пусть H – некоторый набор из $r, r \leq n$, различных признаков вида $\{x_{j_1}, \dots, x_{j_r}\}$. Близость объектов $S' = (a'_1, \dots, a'_n)$ и $S'' = (a''_1, \dots, a''_n)$ из M по набору признаков H будем оценивать величиной $B(S', S'', H)$, принимающей значение 1 , если $a'_{j_t} = a''_{j_t}$ при $t = 1, 2, \dots, r$, и принимающей значение 0 иначе
- **Принципиальная схема построения алгоритмов АВО следующая.**
- В множестве признаков $\{x_1, \dots, x_n\}$ выделяется совокупность различных подмножеств вида $\{x_{j_1}, \dots, x_{j_r}\}$, $r \leq n$, не обязательно одинаковой мощности. В дальнейшем выделенные подмножества называются **опорными множествами** алгоритма, а вся их совокупность обозначается через Ω . Задаются параметры: v_i – параметр, характеризующий представительность (вес), обучающего объекта S_i , $i = 1, 2, \dots, m$, и P_H – параметр, характеризующий представительность (вес) опорного множества H , $H \in \Omega$. Далее проводится **процедура голосования** или **вычисления оценок**. Распознаваемый объект S сравнивается с каждым обучающим объектом S_i , по каждому опорному множеству H .

-

- Пусть S_i – обучающий объект, $S_i \in K$, $K \in \{K_1, \dots, K_l\}$, $H \in \Omega$. Считается, что объект S получает голос за принадлежность классу K , если описания объектов S и S_i совпадают по опорному множеству H (в этом случае $B(S, S_i, H) = 1$). Для каждого класса K , $K \in \{K_1, \dots, K_l\}$, вычисляется оценка принадлежности $\Gamma(S, K)$ объекта S классу K , имеющая вид

- $$\Gamma(S, K) = \frac{1}{N_K} \sum_{S_i \in K} \sum_{H \in \Omega} v_i P_H B(S, S_i, H),$$

- где N_K - число обучающих объектов из класса K .

-

- Объект S относится к тому классу, который имеет наибольшую оценку. Если классов с наибольшей оценкой несколько, то происходит отказ от классификации. Очевидно, что построенный алгоритм не всегда является корректным. Для корректности этого алгоритма требуется выполнение системы линейных неравенств указанного ниже вида.

• Для простоты пусть $l = 2$, $S_i \in K_1$ при $1 \leq i \leq m_1$, $S_i \in K_2$ при $m_1 + 1 \leq i \leq m$, $1 \leq m_1 \leq m - 1$. Тогда система неравенств имеет вид

- $\Gamma(S_1, K_1) > \Gamma(S_1, K_2),$
-
- $\Gamma(S_{m_1}, K_1) > \Gamma(S_{m_1}, K_2),$
-
- $\Gamma(S_{m_1+1}, K_2) > \Gamma(S_{m_1+1}, K_1),$
-
- $\Gamma(S_m, K_2) > \Gamma(S_m, K_1).$
-

• Решение приведённой выше системы неравенств заключается в выборе параметров v_i , $i = 1, 2, \dots, m$, и P_H , $H \in \Omega$. В случае, если система несовместна, находится ее максимальная совместная подсистема и из решения этой подсистемы определяются значения параметров v_i и P_H .

- Другой способ добиться корректности алгоритма – выбрать «хорошую» систему опорных множеств. В частности, выбрать ее так, чтобы для любого обучающего объекта $S' \in K$ было выполнено $G(S', K) > 0$ и для любого обучающего объекта $S'' \notin K$ было выполнено $G(S'', K) = 0$. Это можно сделать следующим образом.
- Пусть H – некоторое опорное множество. Набор признаков H называется *тестом*, если для любых двух обучающих объектов S' и S'' , принадлежащих разным классам, выполнено $B(S', S'', H) = 0$. Другими словами, тест – это набор признаков, по которому различаются любые два объекта из разных классов.
- Пусть совокупность опорных множеств алгоритма есть некоторое множество тестов Ω_T . Очевидно, такой алгоритм является корректным при любых положительных значениях параметров v_i , $i = 1, 2, \dots, m$, и P_H , $H \in \Omega_T$.

- Если набор признаков H_1 – тест, то любой набор признаков H_2 такой, что $H_1 \subset H_2$, также является тестом. При этом если объекты близки по H_2 то они будут близки и по H_1 , если же два объекта близки по набору признаков H_1 , то они не всегда будут близки по H_2 . В этом смысле более короткие тесты обладают большей информативностью, и разумно ограничивать длины тестов или строить тупиковые тесты.
- Набор признаков H называется *тупиковым тестом*, если выполнены следующие два условия 1) H – является тестом; 2) любое собственное подмножество набора H не является тестом. Другими словами, тупиковым тестом является не укорачиваемый набор признаков, по которому любые два обучающих объекта из разных классов отличаются друг от друга.
-
- Оценка $\Gamma(S, K)$ принадлежности распознаваемого объекта S классу K в алгоритме *голосования по тестам* (далее тестовом алгоритме) вычисляется также как и в алгоритме АВО.

- Пусть $\mathbf{H} = \{x_{j_1}, \dots, x_{j_r}\}$ – некоторый набор признаков, $\mathbf{S} = (a_1, \dots, a_n)$ – объект из M . Фрагмент $(a_{j_1}, \dots, a_{j_r})$ описания объекта \mathbf{S} обозначим (\mathbf{S}, \mathbf{H}) .
- Каждый тест \mathbf{H} порождает множество фрагментов описаний объектов вида $(\mathbf{S}_i, \mathbf{H})$, $i = 1, 2, \dots, m$, где \mathbf{S}_i – обучающий объект, причем каждый из этих фрагментов встречается в некотором классе и не встречается в остальных. При распознавании объектов производится голосование по множеству всех таких фрагментов. Рассмотрим простой пример.
- $\{\mathbf{S}_1 = (0, 1, 1, 0), \mathbf{S}_2 = (1, 2, 0, 1), \mathbf{S}_3 = (0, 1, 0, 2)\}$ – класс K_1
- $\{\mathbf{S}_4 = (1, 2, 1, 0), \mathbf{S}_5 = (1, 1, 0, 1), \mathbf{S}_6 = (1, 1, 1, 2)\}$ – класс K_2
- В данном примере наборы признаков $\{x_1, x_2, x_3\}$ и $\{x_1, x_2, x_4\}$ являются тупиковыми тестами. Других тупиковых тестов нет. Если использовать тестовый алгоритм, то объект $\mathbf{S} = (0, 1, 2, 1)$ не будет отнесен ни к одному из классов K_1 и K_2 . Однако $(\mathbf{S}_1, \{x_1, x_2\}) = (\mathbf{S}_3, \{x_1, x_2\}) = (\mathbf{S}, \{x_1, x_2\})$ и нет ни одного объекта из класса K_2 , содержащего такой же фрагмент, что даёт нам основание полагать, что распознаваемый объект \mathbf{S} более близок к классу K_1 .

- Таким образом, если при построении алгоритмов классификации перейти от рассмотрения опорных множеств к анализу фрагментов описаний объектов, то можно строить менее осторожные и при этом корректные процедуры. Примерами таких процедур являются **алгоритмы голосования по представительным наборам** или **алгоритмы типа "Кора"**.
- Фрагмент описания объекта S' из класса K вида (S', H) называется **представительным набором для K** , если для любого обучающего объекта S'' , не принадлежащего классу K , имеет место $B(S', S'', H) = 0$.
- Фрагмент описания объекта S' из класса K вида (S', H) называется **тупиковым представительным набором для K** , если выполнены два условия: 1) для любого обучающего объекта $S'' \notin K$ имеет место $B(S', S'', H) = 0$; 2) для любого набора H' , $H' \subset H$, найдется обучающий объект $S'' \notin K$, для которого $B(S', S'', H') = 1$.

• В классической модели алгоритма голосования по (тупиковым) представительным наборам для каждого класса K строится множество (тупиковых) представительных наборов, обозначаемое далее через $\mathcal{T}(K)$. Распознавание объекта S осуществляется на основе процедуры голосования. В простейшей модификации для оценки принадлежности объекта S классу K используется величина

$$\Gamma(S, K) = \frac{1}{|\mathcal{T}(K)|} = \sum_{(S', H) \in \mathcal{T}(K)} B(S, S', H) .$$

• Очевидно, что все фрагменты описаний обучающих объектов, порожденные некоторым тестом, являются представительными наборами. Очевидно также, что не все представительные наборы, порожденные тупиковым тестом, являются тупиковыми представительными наборами, т.е. в алгоритме голосования по тупиковым представительным наборам строится больше коротких фрагментов описаний объектов, следовательно, он менее осторожный и реже отказывается от распознавания.

- На практике в моделях с представительными наборами отбираются наиболее “весомые” представительные наборы, например те, которые по данному набору признаков встречаются в классе у достаточно большого числа объектов. Рассматриваются также модели с «почти представительными» наборами. В этих моделях информативными для класса K считаются такие фрагменты описаний обучающих объектов, которые по данному набору признаков встречаются достаточно часто в описаниях объектов из класса K и достаточно редко в описаниях объектов из других классов.
- В тестовых моделях для увеличения быстродействия применяются стохастические алгоритмы, в которых построение множества всех тупиковых тестов заменено построением достаточно представительной случайной выборки из него. В данном случае оценки $\Gamma(S, K_j)$, $j = 1, 2, \dots, l$, вычисляются приближённо и оценивается возможная при этом ошибка.

- Рассматриваемые модели могут быть модифицированы на случай вещественнозначной информации путём задания дополнительных параметров, позволяющих устанавливать близость отдельных значений признаков. Для каждого признака x_j , $j = 1, 2, \dots, n$, задаётся некоторый вещественный параметр E_j – точность измерения признака x_j . Опять же в простейшей модификации величина $B(S', S'', H)$ определяется следующим образом. Пусть $S' = (a'_1, \dots, a'_n)$ и $S'' = (a''_1, \dots, a''_n)$, $H = \{x_{j_1}, \dots, x_{j_r}\}$. Тогда $B(S', S'', H) = 1$, если $|a'_{j_t} - a''_{j_t}| \leq E_{j_t}$ при $t = 1, 2, \dots, r$, и $B(S', S'', H) = 0$ иначе.
- Для оценки качества работы распознающего алгоритма A часто используется *процедура скользящего контроля* (leave one out cross validation), которая заключается в следующем. Для каждого i из $\{1, \dots, m\}$ по обучающей выборке $\{S_1, \dots, S_m\} \setminus \{S_i\}$ вычисляются оценки $\Gamma(S_i, K_j)$, $j = 1, 2, \dots, l$. Пусть q – число правильно распознанных объектов S_i . Тогда качество работы алгоритма A оценивается функционалом $\varphi_{\text{ск}}(A) = q/m$.

- Для больших задач описанная процедура скользящего контроля требует существенных вычислительных затрат. В случае, когда A - алгоритм голосования по представительным наборам, для процедуры скользящего контроля существует "быстрый" метод вычисления оценок $\Gamma(S_i, K_j)$, $i = 1, 2, \dots, m, j = 1, 2, \dots, l$. Данный метод позволяет сократить время счета примерно в m раз.
- Для оценки качества работы алгоритма A может быть использован и независимый контроль – набор из t объектов, не вошедших в обучающую выборку, про которые также известно, каким классам они принадлежат. Такой набор объектов называется контрольной выборкой. Пусть q - число правильно распознанных объектов из контрольной выборки. В данном случае качество работы алгоритма A оценивается функционалом $\varphi_{\text{контр}}(A) = q/t$.

- **Замечание.** Среди рассмотренных трёх моделей исторически первыми появились тестовые модели (точнее алгоритмы голосования по множеству тупиковых тестов). Тестовые алгоритмы оказались очень трудоемкими в смысле вычислительных затрат. Теоретические исследования статистических свойств тупиковых тестов показывают, что почти всегда тупиковые тесты имеют примерно одну и ту же длину. Эта длина, вообще говоря, зависит от соотношения между параметрами m и n . Указанное обстоятельство послужило одним из обоснований для построения алгоритмов вычисления оценок, в которых в качестве опорных множеств берутся всевозможные наборы признаков, имеющие фиксированную мощность r . Значение параметра r может задаваться путем некоторого предварительного анализа обучающей выборки. Для этой модели существует формула эффективного вычисления оценки принадлежности $\Gamma(S, K)$ объекта S классу K .

- УПРАЖНЕНИЯ

- 1. Пусть (S', H) – представительный набор класса K . В каком случае при голосовании по (S', H) на скользящем контроле распознаваемый объект S будет отнесен к классу K и в каком случае этот объект будет отнесен к другому классу?

-

- 2. Будем говорить, что признак x_j является фиктивным, если для любого обучающего объекта S' , $S' = (a'_1, \dots, a'_n)$, найдется обучающий объект $S'' = (a''_1, \dots, a''_n)$, принадлежащий тому же классу, что и объект S' и такой, что $a'_j \neq a''_j$, $a'_t = a''_t$ при $t \neq j$, $t \in \{1, 2, \dots, n\}$. Доказать, что фиктивный признак не входит ни в один тупиковый тест.