

ВЫЧИСЛЕНИЕ ВАЖНОСТИ ДИСКРЕТНЫХ ПРИЗНАКОВ (АНАЛИЗ НЕКОТОРЫХ ПОДХОДОВ)

Дана загальна характеристика методів оцінки важливості дискретних ознак. Запропоновано функціонал та загальний алгоритм оцінки важливості, як окремих ознак, так і їх комбінованих груп. Проведена процедура оптимізації запропонованого алгоритму щодо використання обчислювальних ресурсів системи

Ключові слова: завдання розпізнавання, дискретна ознака, оцінка важливості дискретної ознаки (груп ознак)

Дана общая характеристика методов оценки важности дискретных признаков. Предложен функционал и общий алгоритм оценки важности, как отдельных признаков, так и их комбинированных групп. Проведена процедура оптимизации предложенного алгоритма относительно использования вычислительных ресурсов системы

Ключевые слова: задачи распознавания, дискретный признак, оценка важности дискретного признака (групп признаков)

The general characteristic of methods of an estimation of importance of discrete signs is given. It is offered functional and the general algorithm of an estimation of importance, both separate signs, and their combined groups. Procedure of optimization of the offered algorithm concerning use of computing resources of system is spent

Keywords: recognition problems, a discrete sign, an estimation of importance of a discrete sign (groups of signs)

Ф.Г. Ващук

Доктор технических наук, профессор, ректор*

Контактный тел: (0312) 5-15-24

E-mail: vashuk@zakdu.edu.ua

Ю.А. Василенко

Доктор технических наук, профессор

Кафедра информационных управляющих систем и технологий*

Контактный тел: (0312) 2-37-54

E-mail: vasilenko@zakdu.edu.ua

И.Ф. Повхан

Кандидат технических наук, доцент

Кафедра программного обеспечения автоматизированных систем*

Контактный тел: 0685554459

E-mail: comi@zakdu.edu.ua

Л.С. Повхан

Ассистент

Кафедра программного обеспечения автоматизированных систем

Аспирант

Кафедра информационных управляющих систем и технологий*

Контактный тел: 068-734-15-08

E-mail: chajka@zakdu.edu.ua

*Закарпатский государственный университет
ул. Заньковецкой, 87 "Б", г. Ужгород,

Закарпатская область, 88015

1. Введение

Одним из важнейших разделов искусственного интеллекта является распознавание образов. Задача оценки и выбора информативных дискретных признаков некоторых объектов относится к наиболее важным в распознавании. В Украине и за рубежом предложены различные методы уменьшения числа признаков путем выбора некоторого подмножества из исходной совокупности признаков. Однако здесь еще имеется много нерешенных проблем (как вычислительных, так и принципиальных). Например, в [4] указано,

что при наличии трех бинарных признаков "наиболее информативный" из них может не входить в "информативную пару признаков", ниже авторами показана несостоятельность оценки важности отдельных признаков и групп признаков, используемой в тестовых методах и др.

К решению указанной нами задачи известны самые разнообразные подходы. Так, в [6] эта задача решается на основе статистической теории распознавания (предложены приближенные формулы для оценки полезности признаков). В [8] для решения этой задачи используется построение специальной метрики, в ко-

торой расстояние между разделяемыми множествами становится возможно большим. В [7] для определения важности признаков предлагается метод случайного поиска с адаптацией (СПА), а в [1] предложен для этих целей способ, основанный на идее вычисления вероятности значений отдельного признака для каждого класса распознаваемых объектов.

Проблема нахождения хорошей численной характеристики для измерения эффективности данного множества свойств объектов с точки зрения их классификации рассматривалась Бахадуром [5]. “В качестве таких характеристик предлагались различные меры информативности, расстояния и делимости, однако все они не имеют однозначной связи с ошибками классификации”, как отмечено Л.Каналом [4].

В [2, 3] предложен, обоснован и исследован метод оценки важности отдельных дискретных признаков (групп признаков), частично свободный от вышеуказанного недостатка. При этом оценка важности признаков объектов осуществляется непосредственно по характеру компонент разбиения R, количество которых может быть произвольным.

2. Функционал оценки важности дискретных признаков

Оценку важности дискретных признаков p_1, p_2, \dots, p_n по отношению к распознающей функции $f_R(x)$ ($x = p_1 \dots p_n$) наиболее естественно можно определить следующим образом. Пусть M – некоторое множество наборов значений признаков p_1, p_2, \dots, p_n , на каждом из которых известна функция $f_R(x)$. Предположим, что функция $f_R(x)$ принимает значения O_0, O_1, \dots, O_{K-1} . Рассмотрим, например, признак p_1 . Допустим, что этот признак принимает значения $0, 1, \dots, K-1$.

Пусть b_j ($j = 0, 1, \dots, K-1$) – количество всех наборов в множестве M, в которых признак p_1 принимает значение j ; q_j^m – количество всех наборов из M, в которых признак p_1 принимает значение j , а функция $f_R(x)$ принимает значение O_m ($j = 0, 1, \dots, K-1, m = 0, 1, \dots, K-1$); через h обозначим количество всех наборов множества M. Тогда важность $W(p_1)$ признака p_1 по отношению к $f_R(x)$ можно оценить следующей формулой:

$$W(p_1) = \sum_{j=0}^{K-1} \frac{b_j}{h} \cdot \rho_j, \tag{1}$$

где $\rho_j = \max_m (q_j^m / b_j) (0 \leq m \leq K-1)$.

Аналогично оценивается важность остальных признаков p_2, \dots, p_n .

Формула (1) получается из следующих соображений. Величину q_j^m / b_j можно интерпретировать как вероятность того, что функция $f_R(x)$ принимает значение O_m ($0 \leq m \leq K-1$) при условии, что значение признака p_1 равно j ($0 \leq j \leq K-1$). Величина ρ_j представляет собой максимальную из этих вероятностей. Можно еще сказать, что величина ρ_j представляет собой ту информацию, которую можно получить о значении функции $f_R(x)$, зная, что на наборе x значение признака p_1 равно j . Как вытекает из (1), величина $W(p_1)$ определяемая функцией (1), характеризует собой ту информацию, которую можно получить о функ-

ции $f_R(x)$, зная значение признака p_1 на наборе x . Отсюда ясно, что признак p_1 ($1 \leq i \leq n$), для которого эта информация является наибольшей, можно считать самым важным по отношению к функции $f_R(x)$.

Можно показать, что (1) является частным случаем функционала $W(f)$, введенного в [2]. При этом аналогом признака p_1 является обобщенный признак f из [2].

На практике вместо формулы (1) удобно применять следующую формулу:

$$W(p_1) = \frac{1}{h} (\max_m q_0^m + \max_m q_1^m + \dots + \max_m q_{K-1}^m) \tag{2}$$

Формула (1) представляет собой некоторый функционал, определенный на множестве M. Далее (1) будем называть функционалом важности признака p_1 . Аналогично вычисляются функционалы важности признаков p_2, p_3, \dots, p_n .

Правомерность приведенной нами оценки важности дискретных признаков можно убедительно доказать, опираясь на фундаментальные исследования широко известных английских ученых Кендалла и Стюарта (переведены с английского и опубликованы в СССР в 1973 г.). Так в [5] рассмотрены категории, или классы. Приведем для этих данных некоторые сведения из [5], необходимые для дальнейшего изложения затронутой нами темы.

Пусть имеется генеральная совокупность, в которой классификация произведена на основании наличия или отсутствия некоторого признака A. Простейшая задача о взаимозависимости признаков возникает тогда, когда имеется два признака A и B. Если α обозначает отсутствие A и β - отсутствие B, то количества попаданий в четыре возможные подгруппы могут быть в очевидных обозначениях представлены табл. 1.

Эту таблицу 2x2 (иногда называемую 4-х клеточной таблицей) часто удобно записывать в форме табл. 2.

Таблица 1

	B	He B	Сумма
A	AB	Aβ	A
не A	αB	αβ	α
Сумма	B	β	

Таблица 2

a	b	a+b
c	d	c+d
a+c	b+d	n

Если между A и B не существует никакой связи, т.е. если обладание признаком A не связано с обладанием признаком B, то доля индивидов с признаком A среди индивидов, обладающих признаком B, должна быть равна доле индивидов с признаком A среди индивидов, с обладающим признаком B. Таким образом,

по определению, признаки независимы в данной совокупности из наблюдений, если

$$\frac{a}{a+c} = \frac{b}{b+d} = \frac{a+b}{n} \quad (3)$$

Как показано в [5], вероятность правильного оценивания категории А при заданной категории В будет равна

$$\frac{1}{n}(\max(a,c) + \max(b,d)) \quad (4)$$

Аналогично, вероятность правильного оценивания категории В при заданной категории А равна

$$\frac{1}{n}(\max(a,b) + \max(c,d)) \quad (5)$$

Если категоризованные переменные А и В интерпретировать в терминах, использованных нами выше, т. е. под А понимать $f_R(x)$, а под В - p_i ($i=1, n$), то величины a, b, c, d в (4) и (5) будут аналогами $q_0^0, q_0^1, q_1^0, q_1^1$ при $m=k_i=2$ в (1). В этом случае $W(p_i)$ будут представлять собой вероятность правильного оценивания значений $f_R(x)$ при задании признака p_i . Таким образом правомерность предложенной формулы (1) и состоятельность вышеуказанного обоснования $W(p_i)$ не вызывает никаких сомнений.

Используя функционал важности признаков (ФВП), легко реализовать процедуру ранжирования по важности отдельных дискретных признаков p_1, \dots, p_n , характеризующих объекты некоторой заданной обучающей выборки (ОВ). Расположив признаки p_1, \dots, p_n в порядке убывания ФВП для p_i при ($i=1, n$), получим некоторый ранжированный ряд признаков $p_{i1}, p_{i2}, \dots, p_{in}$, где p_{i1} - самый важный признак относительно $f_R(x)$, а p_{in} - наименее важный признак относительно $f_R(p_1, \dots, p_n)$ - функции, задающей разбиение ОВ на образы (классы).

В заключение, отметим некоторые существенные стороны предложенной оценки важности признаков (1).

$$1. \frac{1}{k} \leq W(p_i) \leq 1, \quad (k - \text{число классов ОВ})$$

2. Важность признаков p_1, \dots, p_n в (1) определяется относительно $f_R(x)$ - функции, задающей разбиение множества объектов распознавания, что отсутствует в других известных методах.

3. Численная величина ФВП некоторого признака совпадают с той долей объектов ОВ, которую можно правильно отнести в соответствующие классы на основе учета значений взятого признака, т. е. (1) имеет однозначную связь с ошибками классификации ОВ. Заметим, что свойства 1, 2, 3 непосредственно вытекают из (1).

На основе использования понятия функции принадлежности из теории нечетких множеств были разработаны критерии оценки важности признаков, являющиеся обобщением ФВП.

Оценка важности отдельных дискретных признаков (1) эффективно применяется в методе разветвленного выбора признаков (РВП) [2] для конструирования оптимальных (по памяти и быстродействию)

граф-схем, распознающих наборы дискретных признаков.

3. Оценки важности отдельных групп дискретных признаков

Аналогично приведенной выше оценке важности отдельных признаков (1) можно ввести оценку важности произвольной группы дискретных признаков $p_{i1}, p_{i2}, \dots, p_{iS}$, ($1 \leq i_1, i_2, \dots, i_S \leq n; 2 \leq S \leq n$). Возьмем, например, группу признаков $p_1, p_2, \dots, p_\gamma$, ($2 \leq \gamma \leq n$). Как и выше, М - множество наборов признаков p_1, \dots, p_n , на которых известна функция $f_R(x)$, ($x = p_1, \dots, p_n$), кроме того будем считать, что М является достаточно обширным множеством, т.е. М близко к области определения функции $f_R(x)$.

Пусть $\Delta = t_1, t_2, \dots, t_\gamma$, ($0 \leq t_j \leq k_j - 1; j = 1, 2, \dots, \gamma$) произвольный набор значений признаков. Через B_Δ обозначим количество всех наборов p_1, \dots, p_n из М, для которых выполняется соотношение $p_j = t_j$, $j = 1, 2, \dots, \gamma$ через B_Δ^m - количество наборов p_1, \dots, p_n из М, для которых выполняются соотношения: $p_j = t_j$, при $j = 1, 2, \dots, \gamma$. $f_R(p_1, \dots, p_n) = O_m$ при ($m = 1, 2, \dots, k-1$), через h обозначим количество всех наборов в множестве М, а через Γ - множество всех наборов признаков $p_1, p_2, \dots, p_\gamma$. Тогда важность $W(p_1, \dots, p_\gamma)$ группы признаков определяется по формуле

$$W(p_1, \dots, p_\gamma) = \sum_{\Delta \in \Gamma} \frac{B_\Delta}{h} \sigma_\Delta \quad (6)$$

где $\sigma_\Delta = \max B_\Delta^m / B_\Delta$. Формула (6) обосновывается таким же образом, как и формула (1) (предлагаем читателю проделать это самостоятельно). Можно зафиксировать γ ($1 \leq \gamma \leq n$) и среди всех групп $p_{i1}, p_{i2}, \dots, p_{i\gamma}$ найти самую важную группу признаков. Но при больших γ и n нахождение самой важной группы $p_{i1}, p_{i2}, \dots, p_{i\gamma}$ связано со значительным числом переборов, а именно нужно осуществить:

$$C_n^\gamma = \frac{n!}{\gamma!(n-\gamma)!}$$

переборов. Ниже будут даны более эффективные методы.

Легко видеть, что (6), по аналогии с (1), является частным случаем функционала из [2]. Введенную оценку, по аналогии с ФВП, будем называть функционалом важности группы признаков (ФВПГ).

Отметим, что ФВПГ, как и ФВП, оценивает важность группы признаков относительно $f_R(p_1, \dots, p_n)$ и имеет однозначную связь с ошибками классификации объектов ОВ, т. е. численная величина ФВПГ совпадает с той долей объектов из ОВ, которую можно правильно отнести к образам ОВ, используя только выбранную группу признаков для распознавания объектов (при решающем правиле). Используя ФВПГ, можно легко осуществить процедуру ранжирования подмножеств признаков, что является весьма важной задачей при эксплуатации человеко-машинных распознающих систем диалогового типа. ФВПГ является основой, предложенного автором, аналитического метода нахождения тестов [3], который базируется на следующей теореме.

ТЕОРЕМА. Произвольное подмножество p_1, \dots, p_r , $(1 \leq r \leq n)$ множества признаков p_1, \dots, p_n образует тест ОВ тогда и только тогда, когда $W(p_1, \dots, p_r) = 1$.

Доказательство теоремы приведено в [3]. ФВП и ФВПГ можно использовать для диагностики комбинационных схем при поиске контролируемых и диагностируемых тестов [10]. На основе использования ФВП предложена классификация множества ОВ из дискретных признаков на 2^n типов, характеризующихся определенной сложностью нахождения тупиковых тестов (ТТ), а также табличный метод нахождения ТТ, в котором сокращение перебора достигается за счет выделения некоторого объема арифметических и логических операций, стандартного для всех ОВ, состоящих из объектов с n признаками (двоичными или многозначными). Обобщение ФВП и ФВПГ для случая, когда объекты ОВ есть наборы действительных чисел, а признаки – многозначные предикаты, было дано в [2].

Очевидно, что для ФВПГ имеем: $\frac{1}{K} \leq \text{ФВПГ} \leq 1$, где K - количество образов ОВ.

4. Модифицированный алгоритм оценки важности дискретных признаков

Отметим один важный факт, касающийся вычисления ФВП (ФВПГ) на ПК: для вычисления ФВП совсем не обязательно запоминать все наборы множества M , т.е. вводить в память ЭВМ все объекты ОВ. Величину $W(p_i)$ можно вычислять последовательно при подаче пар $(x, f_R(x))$. Точнее говоря, для вычисления $W(p_i)$ можно предложить следующий алгоритм.

Пусть уже было подано h пар ОВ, и для этих пар вычислены значения b_{ji} и q_{ji}^m , $(j = 0, 1, \dots, K_i - 1; i = 1, 2, \dots, n; m = 0, 1, \dots, K - 1)$; здесь b_{ji} - количество всех наборов из пар ОВ, в которых признак p_i принял значение j , а функция $f_R(x)$ - значение O_m . Тогда при подаче следующей $(h + 1)$ -й пары $(l_1, l_2, \dots, l_n, O_\eta)$, где $0 \leq l_i \leq K_i - 1 (0 \leq \eta \leq K - 1)$, величины h , b_{ji} и q_{ji}^m изменяются следующим образом:

- 1) h изменяется на $h + 1$;
- 2) $b'_{ji} = b_{ji}$, если $j \neq l_i$; $b'_{ji} = b_{ji} + 1$, если $j = l_i$;
- 3) $(q'_{ji})^m = q_{ji}^m$, если $j \neq l_i$ или $m \neq \eta$; $(q'_{ji})^m = q_{ji}^m + 1$, если $j = l_i$ и $m = \eta$.

Таким образом, в памяти ЭВМ можно хранить только величины η , b_{ji} и b_{ji}^m . Вычисление этих величин следует вести до тех пор, пока значения $W(p_i)$, $(i = 1, 2, \dots, n)$ не начнут стабилизироваться вокруг некоторых величин l_i , $(i = 1, 2, \dots, n)$, которые можно принять за оценки важности соответствующих признаков p_i , $(i = 1, 2, \dots, n)$.

Аналогичный алгоритм можно предложить и для ФВПГ. Применение предложенных алгоритмов вычисления важности признаков особенно эффективно в тех условиях, когда имеются определенные ограничения на объем оперативной памяти ЭВМ (например, обработка информации бортовыми ЭВМ).

Уточним связь тестов, их методов распознавания с оценкой важности признаков (1).

Понятие теста, введенное С.В. Яблонским, как известно, является одним из важнейших в теории распознавания. На его основе предложены различные

алгоритмы распознавания образов, а также определяется относительная важность дискретных признаков.

Тесты позволяют не только проанализировать логические связи между признаками, но и вести некоторую меру их важности N_i/N , где N - количество всех тупиковых тестов (ТТ) ОВ, а N_i - количество ТТ ОВ, в которые входит признак p_i (ОВ – обучающая выборка).

Однако отметим, что число вхождений признака p_i в N_i , на наш взгляд, мало характеризует его важность. Действительно, пусть p_i есть ТТ, и, следовательно, число его вхождений в N_i равно 1; в то же время пусть имеется признак p_j , $(j \neq i)$, который, не являясь ТТ, входит во многие ТТ. В этом случае важность p_j по вышеприведенной формуле может намного превышать важность p_i . И в то же время p_i , являясь сам ТТ, безусловно, более важный, чем p_j , так как позволяет различить все объекты ОВ (в отличие от p_j , на основе которого можно различать принадлежность к классам не всех объектов ОВ).

Приведенная мера важности признаков обладает и другими недостатками: огромный объем вычислений при нахождении всех ТТ, эвристичность, необоснованность, применимость в ограниченном числе случаев (когда ОВ содержит малое число объектов и n - невелико).

Численная величина важности признаков $P_i(N_i/N)$ характеризует различие между классами объектов ОВ. Однако, это не всегда выполняется, о чем свидетельствует нижеприведенный авторами пример определения важности признаков ОВ, заданной табл. 3.

Таблица 3

	0	1	2	3	4	5	6	7	8	9	...	15	16	17	...	31
P_1	0	0	0	0	0	0	0	0	0	0		0	1	1		1
P_2	0	0	0	0	0	0	0	0	1	1		1	0	0		1
P_3	0	0	0	0	1	1	1	1	0	0		1	0	0		1
P_4	0	0	1	1	0	0	1	1	0	0		1	0	0		1
P_5	0	1	0	1	0	1	0	1	0	1		1	0	1		1
f_R	0	1	1	0	1	0	0	0	1	0	...	0	1	1	...	1

Здесь уже по самой ОВ непосредственно видно, что наиболее важным признаком для различения объектов двух классов является P_1 (в табл. 3 ОВ состоит из 2-х классов, принадлежность объектов $(P_1 P_2 P_3 P_4 P_5)$ $(p_i \in \{0, 1\})$ к которым задана значением f_R в последней строке: 0 или 1).

Применив общий метод Яблонского, построим для данной ОВ всего один ТТ - $P_1 P_2 P_3 P_4 P_5$. Отсюда, по приведенной выше формуле (N_i/N) , получаем, что важность всех признаков одинакова и равна единице.

Применим формулу (1) для определения важности признака P_1 в ОВ вида табл. 3 в общем случае (для n признаков) $W(p_1) = \frac{1}{2^n} (2^{n-1} + 2^{n-1} - (n-1)) = 1 - \frac{n-1}{2^n}$

При $n = 5$ (табл. 3) имеем $W(p_1) = 1 - \frac{4}{32} \cong 0,9$,

при $n = 10$ имеем $W(p_1) = 1 - \frac{9}{1024} \cong 0,99$.

В то же время важность остальных признаков ОБ табл. 3 - P_2, P_3, P_4, P_5 значительно меньше $W(p_1)$ и равна 0,5 (т.е. минимальной величине важности, которая может быть получена по формуле (1)). Последнее намного ближе к действительности, чем мера N_i/N .

Численная величина $W(p_1)$, как отмечено было выше, совпадает с количеством объектов ОБ (в относительных единицах), которые различаются по принадлежности к классам ОБ на основе учета только значений признака P_1 . Так, в нашем примере (табл. 3) $W(p_1) = \frac{28}{32}$, что означает: при помощи признака P_1 можно правильно распознать 28 объектов из 32-х, допустив ошибку только в случае 4-х объектов (их номера в табл. 3 - 1, 2, 4, 8).

Таким образом, численная величина $W(p_1)$ характеризует собой вполне определенную информацию о степени аппроксимации признаком P_1 характеристической функции $f(p_1, \dots, p_n)$, задающей разбиение, представленное табл. 3.

Последним свойством обладает также и ФВПГ. Так, важность группы признаков P_1, P_2 в ОБ табл. 3 будет равна:

$$W(p_1, p_2) = \frac{1}{32} (\max_m B_{00}^m + \max_m B_{01}^m + \max_m B_{10}^m + \max_m B_{11}^m) = \\ = \frac{1}{32} (5 + 7 + 8 + 8) = \frac{28}{32}$$

(здесь $m = 0, 1$, а $\Gamma = \{00, 01, 10, 11\}$).

Таким образом, при помощи группы признаков P_1, P_2 можно правильно определять принадлежность к соответствующим классам 28-и объектов ОБ табл. 3 (ошибка имеется на объектах с номерами 1, 2, 4, 8).

Из рассмотренного примера вытекает предпочтительность применения меры важности признаков $W(p)$ для практических задач по сравнению с мерой N_i/N : она имеет теоретическое обоснование, просто вычисляется для больших n и M , имеет естественную интерпретацию для группы признаков.

Авторами данной работы реализован ряд алгоритмов рассмотренного здесь подхода (ФВП, ФВПГ) на основе различных алгоритмических языков (Паскаль, С++ и др.).

Полученный программный продукт вошел в качестве автономного дополнения в ПК "Орион", разработанного на базе метода РВП [2].

В заключение отметим, что рассмотренный в работе метод оценки важности дискретных признаков относительно f_R можно применять к обработке (описанию и сжатию) дискретных изображений, определению неисправностей в технических системах, оценке компетентности экспертов (групп экспертов) и ряду других задач, где используются дискретные признаки некоторых объектов (предметов, явлений, ситуаций).

Литература

1. Адасовский Б.И. Метод вычисления информативности многомодальных признаков. – ДАН СССР, 1978, т. 239, № 2. – с. 286 - 289.
2. Василенко Ю.А. Математическое конструирование многоуровневых распознающих систем на основе метода разветвленного выбора признаков: теория, алгоритмы, реализация, применение. Дисс. Докт. техн. наук, Ужгород, 1990.
3. Василенко Ю.А., Шевченко Г.Я. Аналитический метод нахождения тестов. Автоматика, 1979, № 2. – с. 3 - 7.
4. Канал Л. Обзор систем для анализа структуры образов и разработки алгоритмов классификации в режиме диалога. Распознавание образов при помощи цифровых вычислительных машин (пер. с англ.). – М.: Мир, 1974. – с. 124 - 143.
5. Кендал М., Стьюарт А. Статистические выводы и связи (пер. с англ.). – М.: Наука, 1973. – 900 с.
6. Козловский Б.В., Хараузов К.Н. Критерий оценки полезности признаков в задаче распознавания образов. Известия АН СССР. Техническая кибернетика, 1969, № 3. – с. 31 - 36.
7. Лбов Г.С. Выбор эффективной системы зависимых признаков. Вычислительные системы, 1965, вып. 19. – с. 21 - 34.
8. Неймарк Ю.И. К вопросу о выборе признаков при распознавании образов. Известия АН СССР. Техническая кибернетика, 1970, № 1. – с. 41 - 46.