

УДК 519.237

**LRP - ЛОГИЧЕСКОЕ РАСПОЗНАЮЩЕЕ ПРАВИЛО**  
(Описание программы. Решение модельных  
и прикладных задач)

Н.Г. Старцева

В работе дано краткое описание программы LRP, реализующей алгоритм LRP (см. [1]), иллюстрация работы программы на двух модельных примерах, а также сравнение LRP с четырьмя другими известными алгоритмами (на модельных и реальных примерах), которое позволяет убедиться в правильности принципов, заложенных в алгоритме LRP.

§ I. Программа LRP

Программа LRP (оформленная как подпрограмма) предназначена для построения решающего правила в виде логического дерева с заданным количеством вершин, соответствующих обучающей выборке (обучению), последняя (выборка) может быть разбита на  $K$  ( $K \geq 2$ ) образцов. С помощью этого правила распознается контрольная выборка (контроль).

Программа LRP позволяет работать с разнотипными данными и таблицами, в которых есть пропуски. При построении решающего правила LRP одновременно выделяет подсистему информативных признаков. Общие понятия и теоретические исследования так же, как и алгоритм LRP, описаны в [1]. В приложении I дается описание параметров работы LRP и обращение к программе.

Качество работы алгоритма можно определить с помощью приема "скользящий экзамен", реализованного в программе, который состоит в следующем: из обучающей выборки выделяют один элемент и предлагают обучиться на оставшейся части последовательности и классифицировать выделенный элемент, затем выделяют другой элемент (первый возвращается на место) и снова проводят обучение и экзамен на

этом элементе, и так поочередно перебирают все элементы обучающей последовательности. Затем подсчитывается, сколько раз алгоритм ошибался при классификации выделенных элементов. Отношением числа ошибочных классификаций к объему выборки и оценивается качество решающего правила [2].

Рассмотрим режимы работы IRP и параметры, регулирующие качество решающего правила.

В программе реализованы четыре основных режима работы, задаваемые с помощью параметра IR2: IR2=0 - обучение; IR2=1 - обучение и контроль; IR2=2 - скользящий экзамен; IR2=3 - обучение и скользящий экзамен.

В качестве условия деления объектов на каждой вершине дерева рассматриваются простые высказывания или строится гиперплоскость. В зависимости от вида условия деления различаются три режима, устанавливаемые параметром NRE: NRE=1 - режим без гиперплоскости; NRE=2 - режим с гиперплоскостью; NRE=3 - режим только с гиперплоскостью.

В каждой вершине дерева ищется оптимальное в смысле некоторого критерия условия деления объектов. С помощью параметра IR можно выбирать тот или другой критерий деления. При IR=1 в программе IRP реализован только один критерий [2]. При IR=2,3,4,5 можно вставить подпрограмму-функцию, аналогичную реализованной в IEP, результатом которой является другой критерий деления (исследования по выбору новых критериев будут проведены позже).

Алгоритм устроен таким образом, что деление каждый раз происходит из той вершины, у которой значение критерия наименьшее (LD = 0), но в IEP реализован также режим (LD = 1), при котором, если значения критериев для двух вершин отличаются меньше, чем на EPS (обычно EPS = 0), то деление будет происходить из вершины с наименьшей длиной внешнего пути.

Улучшать результаты обучения можно с помощью параметров KALF, KBET, KGAM, KWT, EPS. Параметры KALF =  $\alpha$ , KBET =  $\beta$ , KWT =  $\gamma$ , EPS =  $\epsilon$  приведены в [1]. Через KGAM обозначено минимально возможное количество объектов каждого из двух классов, по которым строится гиперплоскость.

## § 2. Модельный пример и результаты решения

Количество образов в примере равно трем. Объем обучающей выборки - 20 реализаций (по 6 реализаций первого и второго образов

и 8 реализаций третьего образа), объем контрольной выборки - 4 реализации. В обучающей и контрольной выборках есть пропуски. Количество признаков - 6. Первый признак - наименование, четвертый - порядковый, остальные количественные.

Для второго и третьего признаков распределение строится следующим образом: первый образ представляет собой нормальное распределение  $N(\mu_1, \sigma^2)$ , второй -  $N(\mu_2, \sigma^2)$ , третий - смесь двух нормальных распределений  $N(\mu_3, \sigma^2)$  и  $N(\mu_4, \sigma^2)$  с априорными вероятностями  $\frac{1}{4}$  и  $\frac{3}{4}$  соответственно, где  $\mu_1 = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$ ,  $\mu_2 = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$ ,  $\mu_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ ,  $\mu_4 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ ,  $\sigma = 0,62$ .

Первый признак принимает имена "3" или "6" (с равной вероятностью), если при этом второй и третий признаки распределены по нормальному закону  $N(\mu_k, \sigma^2)$ , и имена "1" и "2" (с равной вероятностью) - для всех остальных реализаций. Четвертый признак - порядковый и равен "5" для реализаций, у которых второй и третий признаки распределены по нормальному закону  $N(\mu_3, \sigma^2)$ , и равен "3" или "4" (с равной вероятностью) для всех других реализаций. Пятый и шестой признаки неинформативные (шумовые) и распределены одинаково и независимо для трех образов -  $N(0, 20)$ .

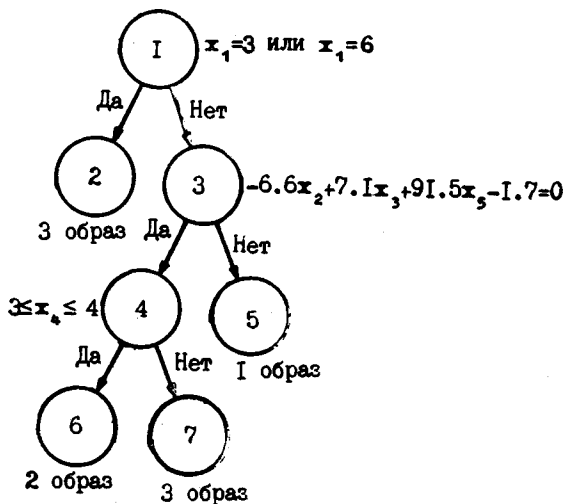


Рис. I

В приложении 2 приводится распечатка результатов работы программы.

На рис. I приведено дерево решений, построенное по результатам обучения.

Построенное решающее правило в виде дерева решений с нулевой ошибкой разбивает объекты обучающей выборки на три класса и является наглядной интерпретацией логических закономерностей данной моделируемой выборки.

§ 3. Результаты сравнения LRP с четырьмя алгоритмами распознавания образов на модельном примере

Количество классов в модельном примере равно двум. Объем обучающей выборки - 100 реализаций (по 50 реализаций каждого образа), объем контрольной выборки - 200 (по 100 реализаций каждого образа). Количество признаков - 20.

Первый признак распределен следующим образом: для первого класса признак попадает в интервалы (0,1), (2,3), (4,5) с вероятностью 0,95; в интервалы (1,2), (3,4), (5,6) с вероятностью 0,05 и распределен внутри каждого интервала равномерно. Для второго класса наоборот. Остальные 19 признаков неинформативны, они распределены одинаково и независимо для обоих классов  $N(0,20)$ .

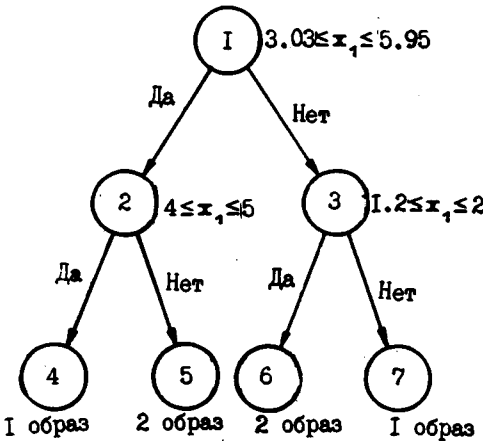


Рис.2

Рассмотрим пять алгоритмов распознавания образов: первый алгоритм использует дискриминантную функцию Фисера [3], второй - использует квадратичную дискриминантную функцию [3], третий основан на непараметрической оценке [4], четвертый - на логических функциях [5], пятый - LRP.

Качество алгоритма определяется оценкой вероятности ошибочной классификации на контроле -  $\bar{P}_\alpha$ , усредненной по числу экспериментов  $n$ ), где  $\alpha$  - номер алгоритма. Усреднение проводилось

\*) Под экспериментом понимается: генерирование обучающей выборки, построение решающего правила, генерирование контрольной выборки и вычисление оценки вероятности ошибки для данного решающего правила и заданного контроля.

по девяти экспериментам. Результаты работы первых четырех алгоритмов на данном модельном примере уже приводились в [6]:  $\bar{P}_1 = 0,45$ ,  $\bar{P}_2 = 0,47$ ,  $\bar{P}_3 = 0,47$ ,  $\bar{P}_4 = 0,18$ . Для алгоритма LBP оценка  $\bar{P}_5 = 0,08$ . Из приведенных выше результатов видно, что алгоритм LBP на данном модельном примере сработал наилучшим образом.

Необходимо отметить, что в решающем правиле распознавания (в узлах дерева) оказываются задействованы не все признаки.

Таким образом, одновременно с построением решающего правила решается задача выбора подсистемы информативных признаков.

На рис.2 приведено дерево решений, построенное по результатам обучения одного из 9 экспериментов, для данного модельного примера.

#### § 4. Результаты сравнения LBP с четырьмя алгоритмами распознавания образов на прикладных задачах

Для эксперимента сравнения были выбраны те же четыре алгоритма, что и в § 3. Качество работы каждого алгоритма определялось оценкой вероятности ошибочной классификации на "скользящем экзамене" -  $\bar{P}_k^c$ . Сравнение пяти алгоритмов проводилось для трех реальных задач: распознавания больных с глубоким наркозом и больных с поверхностным наркозом при операции на сердце (Институт патологии кровообращения ИЗ РСФСР) по анализу концентрации углекислого газа в артериальной крови (первая задача), в капиллярной крови (вторая задача), в венозной крови (третья задача). Данные для третьей задачи представлены в приложении 3. Число признаков в задачах - 6 (все признаки количественные), объем выборки - 14 для больных с глубоким наркозом и 12 - с поверхностным наркозом. В таблице представлены результаты сравнения.

Т а б л и ц а

№ задачи	Вероятность ошибочной классификации				
	$\bar{P}_1^c$	$\bar{P}_2^c$	$\bar{P}_3^c$	$\bar{P}_4^c$	$\bar{P}_5^c$
1	0.39	0.35	0.29	0.61	0.13
2	0.36	0.54	0.29	0.46	0.12
3	0.5	0.54	0.36	0.36	0.11

Необходимо отметить, что при сравнении алгоритмов не учитывались положительные свойства LRP : возможность работать с пропусками в таблицах, с несколькими образами и с разнотипными данными, возможность в вершине дерева рассматривать в качестве условия деления гиперплоскость.

Несмотря на то, что положительные свойства алгоритма в эксперименте сравнения не были учтены, алгоритм LRP показал наилучшее качество решения по сравнению с четырьмя известными алгоритмами. Этот факт объясняется следующим: по сравнению с первым и вторым алгоритмами LRP использует более слабое предположение о виде законов распределения вероятностей; по сравнению с третьим алгоритмом LRP дал лучший результат, так как логические решающие правила являются статистически более устойчивыми [7]; по сравнению с четвертым алгоритмом, также основанным на логических решающих функциях, LRP использует более сложные предикаты (в частности, "двухсторонние интервалы") и имеет более сильный критерий [1].

#### Л и т е р а т у р а

1. ЛБОВ Г.С., СТАРЦЕВА Н.Г. Алгоритмы многоклассового распознавания, основанный на логических решающих функциях. - Настоящий сб., с. 3-10.
2. ВАПНИК В.Н., ЧЕРВОНЕНКИС А.Я. Теория распознавания образов. - М.: Наука, 1974. - 415 с.
3. ФУКУНАГА Е. Введение в статистическую теорию распознавания образов. - М.: Наука, 1979. - 367 с.
4. ЧЕРКАШИН Н.Т. Некоторые непараметрические алгоритмы распознавания образов бальной размерности. - В кн.: Математическая статистика и ее приложение. Томск, 1979, с. 156-162.
5. МАНОХИН А.Н. Методы распознавания образов, основанные на логических решающих функциях. - В кн.: Эмпирическое предсказание и распознавание образов (Вычислительные системы, вып. 67). Новосибирск, 1976, с. 42-53.
6. МАНОХИН А.Н., СТАРЦЕВА Н.Г. Экспериментальное сравнение 4-х алгоритмов распознавания образов. - В кн.: Обнаружение эмпирических закономерностей с помощью ЭВМ (Вычислительные системы, вып. 102). Новосибирск, 1985, с. 127-132.
7. ЛБОВ Г.С. Методы обработки разнотипных экспериментальных данных. - Новосибирск: Наука, 1981. - 160 с.

Поступила в ред.-изд. отд.  
24 мая 1985 года

## Обращение к программе LRF и описание параметров

Обращение к программе: CALL LRF(X,M,N,Y,M1,K,B,IR,KWT,NTP, NRE,LM,KVET,KGAM,KALF,EPS,MOB,LD,IPRINT,NJF,NWX,APF,OSB,NF, F RA2,SI,E,KOB,C,IR2,NOB,ABE,OSE,MR1,RA3,M2,X1).

Входные параметры:

- X - обучающая матрица размерности  $M \times N$ ,
- M - количество реализаций обучающей матрицы,
- N - количество признаков,
- Y - контрольная матрица размерности  $M1 \times N$ ,
- M1 - количество реализаций контрольной матрицы,
- K - количество образов,
- M2 - постоянная, если есть "скользящий экзамен"  $M2 = M - 1$  ; если нет, то  $M2 = 1$ ,
- B - код пропуска (кодируется максимальным положительным значением, не встречающимся в таблице),
- IR - параметр режима, определяет вариант критерия,
- KWT - максимально-возможное количество вершин (целое, нечетное),
- NTP - массив размерности  $N \times 2$ , первый столбец определяет тип признака: 0 - признак наименования, 1 - количественный, 2 - дискретно-количественный, 3 - балльный; второй столбец рабочий,
- NRE - параметр режима работы, определяющий вид выбора условия в вершине,
- LM - максимальное количество имен в одном признаке,
- KVET, KGAM, KALF, EPS, LD - параметры, регулирующие качество решающего правила,
- MOB - массив размерности  $M \times 3$ , первый столбец - вектор принадлежности реализаций к образу, второй столбец - рабочий, третий - выходной - указывает номер вершины, в которую попадает объект обучающей выборки,
- IPRINT - параметр, регулирующий печать: IPRINT = 0 - печати нет, IPRINT = 1 - есть вся печать, IPRINT = 2 - печать только результатов обучения,
- IR2 - параметр режима.

Выходные параметры:

OSS - количество ошибок на обучении,

OSE - количество ошибок на "скользящем экзамене",

ИЛР - массив размерности  $4 \times KWT$ , первый столбец - тип высказывания: 0 - набор имен, 1 - интервал, 2 - набор баллов или дискретов, 3 - гиперплоскость; второй столбец - номер признака, по которому строится высказывание; третий и четвертый - номера признаков для гиперплоскости,

НWX - массив размерности  $K \times KWT$  - определяет количество объектов каждого образа в вершине,

АФР - массив размерности  $4 \times KWT$  пороговых значений, перечень имен для высказываний типа  $\Omega$  (количество имен в наборе не более трех), если один из элементов массива равен В, то для соответствующей вершины это символ конца перечня имен или интервал для высказываний типа 1,2; третий и четвертый элементы строки равны 0 или коэффициенты гиперплоскости для 3,

ИР - массив размерности  $KWT \times 4$ , первый столбец указывает номер вершины, в которую следует идти, если высказывание истинно, второй столбец определяет номер вершины, в которую следует идти, если высказывание ложно, третий столбец указывает номер текущей вершины, четвертый столбец - тип вершины (для конечной вершины значение равно 2),

F - вектор размерности  $KWT$  значений критерия.

МОВ - вектор размерности  $M$  - указывает номер предсказанного образа,

АВЕ - массив условных вероятностей размерности  $M \times K$ .

Рабочие массивы:

RA2 - размерности  $M \times 2$

SI - размерности  $LM \times 2$

B - размерности  $2 \times N$

КОВ - размерности  $K \times 3$

C - размерности  $N \times 2$

MR1 - размерности  $M2 \times 3$

RA3 - размерности  $M2 \times 2$

SI - размерности  $M2 \times N$



Результаты работы программы

Исходные данные:

Количество объектов  $M = 20$   
 Количество признаков  $N = 6$   
 Параметр  $KWBT = 1$   
 Параметр  $KGAM = 2$   
 Параметр  $KALP = 1$   
 Параметр  $KFS = 0,0$   
 Количество классов  $K = 3$   
 Количество ветвей  $KWT = 7$   
 Номер критерия  $IR = 1$   
 Длина внешнего пути есть  $LD = 1$   
 Пропуск  $B = 9,0$

Распределение объектов по образам в каждой вершине:

Вершина номер 1	6,	6,	8
Вершина номер 2	0,	0,	6
Вершина номер 3	6,	6,	2
Вершина номер 4	0,	6,	2
Вершина номер 5	6,	0,	0
Вершина номер 6	0,	6,	0
Вершина номер 7	0,	0,	2

Результат распознавания:

Количество объектов на контроле  $M1 = 4$   
 Объект номер 1 принадлежит образу 3  
 Условные вероятности принадлежности к образам

0.0	0.0	1.000
-----	-----	-------

Объект номер 2 принадлежит образу 3

Условные вероятности

0,0            0.0            1.000

Объект номер 3 принадлежит образу 1

Условные вероятности

0.4286            0.4286            1.1429

Объект номер 4 принадлежит образу 3

Условные вероятности

0.0            0.0            1.000

## Результат обучения

K5	Р е ш а ю щ е е п р а в и л о					
	NF(K5,3)	NJF(1,K5)	NJF(2,K5)	NJF(3,K5)	NJF(4,K5)	APF(1,K5)
1	1	0	1	0	0	3.0000
2	2	0	0	0	0	0.0
3	3	3	2	3	5	-6.6261
4	6	0	0	0	0	0.0000
5	7	0	0	0	0	0.0
6	4	2	4	0	0	3.00
7	5	0	0	0	0	0.0

K5	Р е ш а ю щ е е п р а в и л о					
	APF(2,K5)	APF(3,K5)	APF(4,K5)	NF(K5,1)	NF(K5,2)	NF(K5,4)
1	6.0000	9.0	0.0	2	3	1
2	0.0	0.0	0.0	0	0	2
3	7.0922	-1.4888	-1.6850	4	5	1
4	0.0000	0.0	0.0	0	0	2
5	0.0	0.0	0.0	0	0	2
6	4.00	0.0	0.0	6	7	1
7	0.0	0.0	0.0	0	0	2

Здесь: K5 - номер, NF(K5,3) - номер исходящей вершины, NF(K5,1) - номер вершины, в которую надо идти, если высказывание истинно, NF(K5,2) - номер вершины, в которую надо идти, если высказывание ложно, NF(K5,4) - указывает конечность вершины, NJF(1,K5) - индекс высказывания, NJF(J,K5), J = 2, 4, - номера признаков, по которым строится высказывание, APF(J,K5), J = 1, 4, - пороги. Количество ошибок при обучении OSS = 0.0.

Т а б л и ц а 2

№ объекта	№ признака					
	1	2	3	4	5	6
1	5.0	6.0	5.6	6.0	5.8	0.0
2	0.3	0.3	5.7	6.0	6.1	3.7
3	0.5	5.6	5.0	1.8	14.0	6.8
4	0.0	12.1	17.5	0.0	17.0	13.0
5	3.9	0.0	4.7	0.0	17.2	9.1
6	7.4	5.4	7.5	0.0	22.2	18.0
7	6.8	10.0	0.0	0.0	0.0	8.5
8	8.0	11.2	10.1	14.8	13.3	4.4
9	4.8	6.3	0.0	7.8	18.8	1.6
10	0.0	0.0	3.8	3.9	7.6	3.1
11	11.0	0.0	0.0	0.0	22.1	11.7
12	8.0	11.6	0.0	0.0	10.0	7.5
13	7.5	0.0	3.5	5.0	0.0	5.0
14	0.0	0.1	8.2	0.0	8.2	4.8
15	0.0	5.8	5.6	6.0	19.9	12.2
16	0.0	2.2	0.0	11.8	20.3	10.9
17	2.2	0.5	3.0	0.0	15.0	10.0
18	0.0	2.9	0.0	8.0	16.4	12.8
19	0.5	0.0	0.0	5.6	14.9	7.5
20	0.0	6.0	10.0	0.0	14.8	13.2
21	0.0	0.0	0.0	0.0	15.0	10.8
22	6.1	0.1	0.0	7.5	4.8	2.8
23	0.0	0.0	7.6	9.0	18.2	11.5
24	8.8	0.0	11.1	0.0	22.7	9.3
25	0.0	7.2	7.8	0.0	0.0	12.8
26	9.5	11.5	0.0	0.0	13.6	5.5
27	4.4	0.0	0.0	9.8	15.3	7.7
28	0.0	6.2	9.5	6.2	18.0	9.5