

ОБ ОДНОМ ПОНЯТИИ СЛОЖНОСТИ СТРАТЕГИИ
ПРИРОДЫ В РАСПОЗНАВАНИИ ОБРАЗОВ

Г.С.Лбов, Н.Г.Старцева

В в е д е н и е

Из теоретических исследований большого числа авторов следует, что чем более сложные функциональные зависимости используются для построения решающих правил распознавания, тем больше привлекается число признаков и чем меньше объем выборки, тем больше вероятность получения "плохого" правила, сильно отличающегося от оптимального. Вопрос о соотношении сложности класса решающих правил и объема выборки является наиболее важным и трудным в общих теоретических исследованиях, связанных с принятием решений на основе ограниченной выборки. Эта проблема получила название проблемы устойчивости статистических решений.

В настоящее время наметились в основном два пути решения указанной проблемы. В рамках первого пути [1] вводится понятие меры сложности класса решающих правил. Для оценки качества решающего правила предлагается критерий, который одновременно учитывает оценку вероятности ошибки, полученной на обучающей выборке, меру сложности класса решающих правил и объем выборки. Однако рекомендуемые при таком подходе для получения устойчивых решений объемы выборок являются достаточно большими (значительно превышающими объемы выборок в прикладных задачах). Отсюда следует, что либо практически во всех прикладных задачах мы в принципе не можем получить устойчивые решения, либо предложенная теоретическая модель неадекватно отражает действительность. С нашей точки зрения, эта неадекватность заключается в том, что используемая мера сложности класса решающих правил ориентируется на самый худший случай расположения выборки в многомерном пространстве, вероятность которого для прикладных задач мала или равна нулю.

В рамках другого подхода [2] решение указанной проблемы рассматривается в предположении, что распределение вероятностей в пространстве признаков (стратегия природы) выбирается из некоторого ограниченного класса распределений (например, нормального). В этом случае для получения решающего правила, близкого к оптимальному, для наиболее часто используемых на практике алгоритмов распознавания требуется небольшой объем выборки. Это является достоинством данного подхода. Однако при его использовании возникает вопрос о соответствии предполагаемого класса распределений действительному [3,4].

В данной работе сделана попытка обойти трудности, связанные с применением обоих подходов. Идея предлагаемого подхода заключается в следующем. Вводится мера сложности класса распределений вероятностей. Эта мера позволяет получить упорядоченную по сложности последовательность вложенных классов распределений вплоть до класса, содержащего любое произвольное распределение, и исследовать известные алгоритмы на устойчивость в зависимости от меры сложности распределений. Указанное исследование проведено пока для одномерного случая.

§1. Основные определения

Проблему построения решающего правила в распознавании образов можно рассматривать со статистических позиций [5]. Пусть имеется несколько гипотез, каждой из которых соответствует свое распределение вероятностей для наблюдений. Необходимо принять одну из этих гипотез, отвергнув остальные. Не нарушая общности, в дальнейшем будем рассматривать случай двух гипотез. При построении решающего правила желательно минимизировать вероятность ошибочной классификации.

Пусть A - множество объектов из $\pi_1 \cup \pi_2$, причем наблюдаемый объект $a \in A$ относится либо к генеральной совокупности π_1 , либо к генеральной совокупности π_2 . Будем говорить, что объект a относится к i -му классу, если $a \in \pi_i$ ($i = 1, 2$). Функции $X_1, X_2, \dots, \dots, X_j, \dots, X_n, X_{n+1}$ такие, что $X_j: A \rightarrow R$, где R - множество действительных чисел, называются признаками. Признак X_{n+1} - целевой и принимает два значения: $X_{n+1} = 1$ или $X_{n+1} = 2$. В дальнейшем будем рассматривать одномерное признаковое пространство ($n=1$). Конкретным значением признака X для объекта a будет величина $x = X(a)$, $x \in R$.

ОПРЕДЕЛЕНИЕ 1. Решающее правило Ψ есть отображение, которое ставит в соответствие точке x значение $i \in \{1, 2\}$.

Согласно решающему правилу Ψ объект будет отнесен либо к генеральной совокупности π_1 , либо к π_2 .

ОПРЕДЕЛЕНИЕ 2. Под стратегией природы s , соответствующей генеральным совокупностям π_1 и π_2 , будем понимать пару $s = \{p(1, x), p(2, x)\}$, где $p(i, x) = q_i p_i(x)$, $i = 1, 2$, q_i - априорная вероятность появления объекта из π_i , $p_i(x)$ - условная плотность распределения вероятности, соответствующая генеральной совокупности π_i .

Не нарушая общности, в дальнейшем будем рассматривать случай равных априорных вероятностей ($q_1 = q_2 = \frac{1}{2}$).

Стратегия природы s задает вероятностную характеристику генеральных совокупностей π_1 и π_2 .

§2. Понятие сложности стратегии природы

В качестве решающего правила Ψ рассмотрим оптимальное байесовское решающее правило Ψ_0 :

$$\Psi_0(x) = \begin{cases} 1 & \text{при } \frac{p_1(x)}{p_2(x)} \geq \frac{q_2}{q_1}, \\ 2 & \text{при } \frac{p_1(x)}{p_2(x)} < \frac{q_2}{q_1}. \end{cases}$$

Вероятность ошибочной классификации для байесовского решающего правила является минимальной и определяется в виде:

$$\mathcal{P}_0 = q_1 \int_{x \in E_2} p_1(x) dx + q_2 \int_{x \in E_1} p_2(x) dx,$$

где $E_1 = \{x | \Psi_0(x) = 1\}$, $E_2 = \{x | \Psi_0(x) = 2\}$, $E_1 \cup E_2 = R$, $E_1 \cap E_2 = \emptyset$. Фиксируя стратегию природы s , мы однозначно задаем \mathcal{P}_0 . Обратное утверждение, вообще говоря, неверно, т.е. двум различным стратегиям s_1 и s_2 может соответствовать одинаковая по величине вероятность ошибочной классификации \mathcal{P}_0 .

Введем класс $S_{\mathcal{P}_0}$ всех стратегий при фиксированной вероятности ошибочной классификации. Тогда полный класс стратегий природы $S = \bigcup_{\mathcal{P}_0} S_{\mathcal{P}_0}$.

Байесовское решающее правило, построенное для некоторой произвольной стратегии $s \in C_{\mathcal{P}_0}$, разбивает диапазон изменения значений признака X на $L+1$ интервал L границами (L может быть сколь угодно большим). При этом всем значениям признака X из одного интервала приписывается некоторое одинаковое решение: либо $\Psi(x) = 1$, либо $\Psi(x) = 2$.

ОПРЕДЕЛЕНИЕ 3. Под сложностью произвольной стратегии природы $s \in C_{\mathcal{P}_0}$ будем понимать некоторую функцию $l = l(\gamma, \mathcal{P}_0)$, где l — минимальное количество границ из общего числа границ, соответствующих байесовскому решающему правилу. Вероятность ошибочной классификации \mathcal{P}_γ , полученная при l границах, отличается от \mathcal{P}_0 меньше, чем на γ .

Очевидно, что такое определение сложности стратегии природы возможно при фиксированных γ и \mathcal{P}_0 .

УТВЕРЖДЕНИЕ I. Для двух произвольных стратегий природы s_1 и $s_2 \in C_{\mathcal{P}_0}$ порядок сложности $l(\gamma, \mathcal{P}_0)$ с изменением γ может меняться.

ДОКАЗАТЕЛЬСТВО. Пусть s_1 и $s_2 \in C_{\mathcal{P}_0}$ и сложность стратегии s_1 равна $l_1(\gamma_1, \mathcal{P}_0) = l_1$. Пусть s_1 такова, что при $\gamma_2 = 2\gamma_1$ сложность ее не меняется. Такую стратегию s_1 всегда можно найти в $C_{\mathcal{P}_0}$. Пусть s_2 имеет сложность $l_2(\gamma_1, \mathcal{P}_0) = l_1 + 1$. При $\gamma_2 = 2\gamma_1$ сложность стратегии s_2 уменьшается, т.е. $l_2(\gamma_2, \mathcal{P}_0) < l_1 + 1$. Такая стратегия в классе $C_{\mathcal{P}_0}$ всегда найдется. Таким образом, при $\gamma = \gamma_1$: $l_1(\gamma_1, \mathcal{P}_0) = l_1 < l_1 + 1 = l_2(\gamma_1, \mathcal{P}_0)$, но $l_1(\gamma_2, \mathcal{P}_0) = l_1 \geq l_2(\gamma_2, \mathcal{P}_0)$, так как $l_2(\gamma_2, \mathcal{P}_0) < l_1$, что и требовалось доказать.

Согласно утверждению I упорядочивание стратегий природы s имеет место только при фиксированных \mathcal{P}_0 и γ .

УТВЕРЖДЕНИЕ 2. Для стратегии природы $s \in C_{\mathcal{P}_0}$ с ростом γ сложность стратегии $l(\gamma, \mathcal{P}_0)$ есть функция невозрастающая.

Доказательство следует из определения 3 сложности стратегии природы, т.е. если $\gamma_1 > \gamma_2$, то $l(\gamma_1, \mathcal{P}_0) \leq l(\gamma_2, \mathcal{P}_0)$.

ОПРЕДЕЛЕНИЕ 4. Произвольное решающее правило Ψ называется квазиоптимальным порядка ε , если для стратегии $s \in C_{\mathcal{P}_0}$ с заданной сложностью $l(\gamma, \mathcal{P}_0)$ величина

$\mathcal{P}' - \mathcal{P}_0 \leq \epsilon$, где \mathcal{P}_0 - вероятность ошибочной классификации байесовского решающего правила, \mathcal{P}' - вероятность ошибочной классификации квазиоптимального решающего правила.

Очевидно, что при $\epsilon = 0,5$ все решающие правила являются квазиоптимальными.

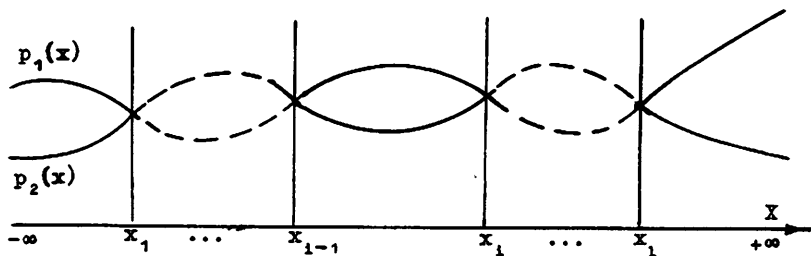
ОПРЕДЕЛЕНИЕ 5. Под алгоритмом распознавания понимается процедура выбора решающего правила из некоторого фиксированного класса правил.

ОПРЕДЕЛЕНИЕ 6. Будем говорить, что алгоритм S решает задачу распознавания образов сложности $l(\gamma, \mathcal{P}_0)$ с точностью ϵ , если любое решающее правило Ψ_S , порожденное алгоритмом S , является квазиоптимальным порядка ϵ для некоторой $c \in \mathcal{C}_{\mathcal{P}_0}$.

Фиксируя \mathcal{P}_0 и γ , можно дать рекомендации пользователю, какой сложности стратегии природы и с какой точностью он может распознавать исследуемый алгоритм S .

§3. Усредненная стратегия природы

Рассмотрим произвольную стратегию $c \in \mathcal{C}_{\mathcal{P}_0}$ с фиксированной сложностью $l(\gamma=0, \mathcal{P}_0)=1$ (см. рисунок).



Здесь x_i - i -я граница, соответствующая байесовскому решающему правилу ($i = 1, 1$).

Для стратегии $c \in \mathcal{C}_{\mathcal{P}_0}$ вероятность ошибочной классификации байесовского решающего правила

$$\mathcal{P}_0 = \frac{1}{2} \sum_{i=1}^{l+1} \min(p_1^i, p_2^i),$$

где

$$p_1^i = \int_{x_{i-1}}^{x_i} p_1(x) dx, \quad p_2^i = \int_{x_{i-1}}^{x_i} p_2(x) dx, \quad i = \overline{1, l+1} \quad (x_0 = -\infty, x_{l+1} = +\infty).$$

Здесь

$$\sum_{i=1}^{l+1} p_j^i = 1; \quad j = 1, 2, \quad i = \overline{1, l+1}, \quad 0 \leq p_j^i \leq 1.$$

УТВЕРЖДЕНИЕ 3. Значение вероятности ошибочной классификации байесовского решающего правила \mathcal{P}_0 одна и та же для любых плотностей распределения $p_1(x), p_2(x)$ внутри интервала (x_{i-1}, x_i) при фиксированных p_1^i и $p_2^i, i = \overline{1, l+1}$.

В дальнейшем будем считать (следует из утверждения 3), что стратегии природы, при которых байесовская процедура дает 1 граници и \mathcal{P}_0 при фиксированных p_1^i и $p_2^i, i = \overline{1, l+1}$, являются неразличимыми.

Введем равномерную вероятностную меру на $\mathcal{C}_{\mathcal{P}_0}$.

ОПРЕДЕЛЕНИЕ 7. Под "усредненной стратегией природы" будем понимать следующую стратегию: 1, \mathcal{P}_0 фиксировано, $\gamma = 0$; если l четно, то

$$p_1^i = \frac{2(1 - \mathcal{P}_0)}{l+2}; \quad p_2^i = \frac{\mathcal{P}_0}{l+2} \quad \text{для } i = 1, 3, \dots, l+1,$$

$$p_1^i = \frac{\mathcal{P}_0}{1}; \quad p_2^i = \frac{2(1 - \mathcal{P}_0)}{1} \quad \text{для } i = 2, 4, \dots, l;$$

если l нечетно, то

$$p_1^i = \frac{2(1 - \mathcal{P}_0)}{l+1}; \quad p_2^i = \frac{\mathcal{P}_0}{l+1} \quad \text{для } i = 2, 4, 6, \dots, l;$$

$$p_1^i = \frac{\mathcal{P}_0}{l+1}; \quad p_2^i = \frac{2(1-\mathcal{P}_0)}{l+1} \quad \text{для } i = 1, 3, \dots, l+1.$$

Внутри каждого интервала $[x_{l-1}, x_l]$ плотность распределения по каждому из образов равномерна.

Для "усредненной" стратегии будем рассматривать область определения признака X на интервале от 0 до 1. Любую стратегию можно отобразить на этот интервал.

В дальнейшем из всего множества стратегий $C_{\mathcal{P}_0}$ будем рассматривать только "усредненные" стратегии.

§4. Экспериментальное сравнение пяти алгоритмов распознавания образов

Пусть величина \mathcal{P}_0 принимает значения из множества $\{0; 0,05; 0,1; 0,15\}$. При фиксированном \mathcal{P}_0 рассмотрим "усредненные" стратегии сложности $l = 1, 2, 3, 4, 9$. Качество алгоритма S будем определять, как это широко принято [2,6,7], математическим ожиданием вероятности ошибочной классификации MP решающих правил, порожденных этим алгоритмом на всевозможных случайных выборках данного объема.

ОПРЕДЕЛЕНИЕ 8. Под машинным экспериментом понимается последовательность таких процедур, как генерирование обучающей выборки, построение решающего правила, порожденного алгоритмом S , генерирование контрольной выборки и вычисление оценки вероятности ошибки для заданного контроля.

Для каждой "усредненной" стратегии будем рассматривать оценку вероятности ошибочной классификации на контроле \bar{P}_S , усредненную по числу экспериментов (здесь 9 экспериментов).

Тогда определение 6 можно переписать следующим образом.

ОПРЕДЕЛЕНИЕ 9. Некоторый произвольный алгоритм S распознает с точностью ϵ "усредненную" стратегию $s \in C_{\mathcal{P}_0}$ с заданной сложностью l (решает задачу распознавания образов сложности l), если $\bar{P}_S - \mathcal{P}_0 \leq \epsilon$.

Такой алгоритм S в дальнейшем будем называть квазиоптимальным порядка ϵ для усредненной стратегии $s \in C_{\mathcal{P}_0}$ со сложностью l .

Количество классов в эксперименте равно двум. Объем обучающей выборки – 100 реализаций (по 50 реализации каждого образа), объем контрольной выборки – 200 (по 100 реализации каждого образа). Решаются две задачи: 1) рассматриваются "усредненные" стратегии, заданные на одном признаке; 2) рассматриваются 20 признаков: первый признак – тот же, а остальные 19 признаков неинформативны и распределены одинаково и независимо по нормальному закону для обоих классов: $N(0,20)$ (стратегия с "шумами").

Были рассмотрены пять алгоритмов распознавания образов: первый ("Фишер") использует дискриминантную функцию Фишера [6], второй ("Квадрат") использует квадратичную дискриминантную функцию [6], третий ("Парзен") основан на непараметрической оценке [8], четвертый (DW13 [9]) и пятый (LRP[10]) основаны на логических решающих функциях*).

Необходимо отметить, что при сравнении алгоритмов не учитывались положительные свойства LRP и DW13 – возможность работать с разнотипными данными (количественными, порядковыми, номинальными, булевыми), а также для LRP: возможность работать с пропусками в таблицах, с несколькими образами, возможность в вершине дерева рассматривать гиперплоскость.

§5. Результаты экспериментального сравнения пяти алгоритмов распознавания образов

В результате машинного эксперимента были получены результаты, представленные в табл. 1 и 2:

1) в табл.1 при различных значениях \mathcal{P}_0 и 1 для "усредненных" стратегий указаны значения ϵ , при которых алгоритмы становятся квазиоптимальными;

2) в табл.2 представлены аналогичные результаты табл.1, но для "усредненных" стратегий с "шумами".

В дальнейшем для обобщения результатов будем рассматривать квазиоптимальность порядка $\epsilon = 0,05$.

Из этих таблиц можно сделать следующие выводы:

1) с увеличением сложности стратегии при фиксированной \mathcal{P}_0 порядок квазиоптимальности не уменьшается для всех алгоритмов. Исключения из этого правила указаны ниже (см. пп. 8,9,10);

*) Ранее в [7] было проведено сравнение 13 наиболее употребительных на практике алгоритмов распознавания. В результате были отобраны в качестве наилучших три первых из перечисленных выше алгоритмов.

2) быть квазиоптимальными порядка $\epsilon = 0,05$ в условиях "шума" могут только алгоритм LRP для всех значений \mathcal{P}_0 и $l = 1, 2, 3, 4$ и алгоритм DW13 для $\mathcal{P}_0 = 0,05$ и $l = 1, 2, 3, 4$. Это и не удивительно, так как алгоритмы, основанные на логических решающих функциях, строят решающее правило с одновременным выбором информативных признаков;

3) для алгоритмов "Фишер", "Квадрат" и "Парзен" предварительно необходимо делать отбор информативных признаков, а далее проводить сравнение;

4) для стратегий сложности $l = 1$ при любом значении \mathcal{P}_0 все алгоритмы являются квазиоптимальными порядка $\epsilon = 0,05$;

5) для стратегий сложности $l = 2$ при любом значении \mathcal{P}_0 все алгоритмы, кроме "Фишера", являются квазиоптимальными порядка $\epsilon = 0,05$;

6) для стратегий сложности $l = 3$ при любом значении \mathcal{P}_0 алгоритмы LRP, DW13, "Парзен" являются квазиоптимальными порядка $\epsilon = 0,05$;

7) самые сложные ($l = 9$) стратегии распознает с точностью $\epsilon = 0,05$ только LRP при всех значениях \mathcal{P}_0 , кроме $\mathcal{P}_0 = 0,15$, а при $\mathcal{P}_0 = 0; 0,05$ еще и "Парзен";

8) для квадратичного решающего правила при всех значениях \mathcal{P}_0 порядок квазиоптимальности уменьшается до величины 1, кратной двум, в случае "нешумовых" стратегий (первая задача);

9) в условиях "нешумовых" стратегий алгоритм "Фишер" проводит гиперплоскость посередине интервала $[0, 1]$. Это можно легко увидеть, расписав уравнение гиперплоскости [6]. Поэтому с увеличением сложности (при $l = 9$) порядок квазиоптимальности резко падает. Для $l = 2, 3, 4$ этот алгоритм имеет $\bar{F}_1 \approx 0,5$, но при $l \rightarrow M$, где M - объем обучающей выборки, снова $\bar{F}_1 \approx 0,5$;

10) в условиях "шумовых" стратегий гиперплоскость для алгоритма "Фишер" проходит не обязательно посередине, этим можно объяснить падение квазиоптимальности при $l = 3$;

11) наилучшее качество решения LRP по сравнению с DW13 объясняется тем, что LRP использует более сложные предикаты (в частности, "двухсторонние интервалы") и имеет более "сильный" критерий [10].

Т а б л и ц а I

	Мял алгоритма	l = 1	l = 2	l = 3	l = 4	l = 9
$\varphi_0 = 0$	"Фингер"	0.02	0.5	0.5	0.5	0.4
	"Квадрат"	0.02	0.06	0.49	0.26	0.40
	"Парзен"	0.02	0.03	0.03	0.03	0.04
	DW13 LRP	0.01 0.004	0.02 0.004	0.04 0.01	0.04 0.02	0.25 0.035
$\varphi_0 = 0,05$	"Фингер"	0.02	0.42	0.45	0.45	0.36
	"Квадрат"	0.02	0.05	0.44	0.24	0.4
	"Парзен"	0.01	0.04	0.04	0.04	0.06
	DW13 LRP	0.02 0.006	0.02 0.01	0.04 0.03	0.04 0.03	0.18 0.04
$\varphi_0 = 0,1$	"Фингер"	0.01	0.39	0.4	0.4	0.32
	"Квадрат"	0.01	0.06	0.38	0.2	0.33
	"Парзен"	0.01	0.03	0.03	0.03	0.07
	DW13 LRP	0.03 0.003	0.04 0.007	0.05 0.02	0.05 0.02	0.16 0.04
$\varphi_0 = 0,15$	"Фингер"	0.008	0.34	0.35	0.33	0.30
	"Квадрат"	0.01	0.05	0.34	0.2	0.34
	"Парзен"	0.008	0.03	0.04	0.04	0.21
	DW13 LRP	0.04 0.009	0.04 0.02	0.05 0.02	0.05 0.02	0.19 0.17

Т а б л и ц а 2

	Имя алгоритма	1 = 1	1 = 2	1 = 3	1 = 4	1 = 5
$\varphi_0 = 0$	"Фигер"	0,06	0,5	0,41	0,47	0,5
	"Квадрат"	0,13	0,32	0,43	0,47	0,45
	"Парзен"	0,13	0,27	0,04	0,44	0,5
	DW13	0,003	0,02	0,03	0,05	0,44
	LRP	0,002	0,007	0,008	0,008	0,3
$\varphi_0 = 0,05$	"Фигер"	0,08	0,43	0,35	0,42	0,43
	"Квадрат"	0,28	0,35	0,43	0,44	0,45
	"Парзен"	0,22	0,32	0,38	0,4	0,43
	DW13	0,03	0,03	0,05	0,05	0,37
	LRP	0,003	0,01	0,03	0,03	0,35
$\varphi_0 = 0,1$	"Фигер"	0,08	0,36	0,29	0,37	0,39
	"Квадрат"	0,28	0,32	0,38	0,39	0,39
	"Парзен"	0,23	0,29	0,32	0,36	0,38
	DW13	0,08	0,08	0,08	0,15	0,38
	LRP	0,003	0,006	0,03	0,03	0,31
$\varphi_0 = 0,15$	"Фигер"	0,10	0,32	0,33	0,3	0,35
	"Квадрат"	0,26	0,24	0,28	0,34	0,35
	"Парзен"	0,26	0,23	0,33	0,3	0,35
	DW13	0,1	0,1	0,13	0,13	0,29
	LRP	0,04	0,04	0,04	0,04	0,22

Л и т е р а т у р а

1. ВАПНИК В.Н., ЧЕРВОНЕНКИС А.Я. Теория распознавания образов. М.: Наука, 1974. - 415 с.
2. РАУДИС Ш. Ограниченность выборки в задачах классификации. -В кн.: Статистические проблемы управления, вып. 18, Вильнюс, 1976.
3. HUGHES G.P. On the mean accuracy of statistical pattern recognizers. -IEEE Trans. Inform. Theory, 1968, v. IT-14, p. 55-63.
4. ДУДА Р., ХАРТ П. Распознавание образов и анализ сцен. -М.: Мир, 1976. - 559 с.
5. АНДЕРСЕН Г. Введение в многомерный статистический анализ. -М.: ФМ, 1963. - 472 с.
6. ФУКУНАГА. Введение в статистическую теорию распознавания образов. -М.: Наука, 1979. - 367 с.
7. РАУДИС Ш., ПИЖАЛИС В., ЮШКЕВИЧЮС К. Экспериментальное сравнение тринадцати алгоритмов классификации. -В кн.: Статистические проблемы управления, вып. 13, Вильнюс, 1975.
8. ЧЕРКАШИН Н.Т. Некоторые непараметрические алгоритмы распознавания образов большой размерности. -В кн.: Математическая статистика и ее приложение. Томск, 1979, с. 156-162.
9. МАНОХИН А.Н. Методы распознавания образов, основанные на логических решающих функциях. -В кн.: Эмпирическое предсказание и распознавание образов (Вычислительные системы, вып. 67). Новосиби- бирск, 1976, с. 42-53.
10. ЛБОВ Г.С., СТАРЦЕВА Н.Г. Алгоритм многоклассового распознавания, основанный на логических решающих функциях. -В кн.: Методы анализа данных (Вычислительные системы, вып. III). Новосиби- бирск, 1985, с. 3-II.

Поступила в ред.-изд.отд.
28 мая 1986 года