

ВЫЧИСЛИТЕЛЬНЫЕ СИСТЕМЫ

Сборник трудов

Института математики СО АН СССР

1966 г.

Выпуск 22

ПОЛУЧЕНИЕ ДОСТАТОЧНЫХ ХАРАКТЕРИСТИК

$\rho_i(x)$ ПРИ РАСПОЗНАВАНИИ ОБРАЗОВ

В.М. Курилов

Статистическая постановка задачи распознавания образов не нова [1,2,3]. Основным недостатком существующей постановки является требование точного знания плотностей распределения вероятностей образов и отсутствие метода получения достаточных оценок этих плотностей в процессе обучения по конечной выборке. Кроме того, в классической постановке не дается ответа на вопрос, какова должна быть выборка обучающей последовательности.

Это заставляет вернуться к статистической трактовке задачи распознавания образов.

Введем некоторые понятия и обозначения, которые понадобятся в дальнейшем.

Под образом будем понимать некоторое подмножество из множества отражений предметов или явлений материального мира, выделенное воспринимающей системой в пространстве измеряемых ею параметров в замкнутую односвязную область.

Каждое отражение предмета (или явления) из этого подмножества будем называть реализией образа. Реализация характеризуется n -мерным вектором в пространстве измеряемых системой параметров. Значение этого вектора есть случайная величина. Тогда задачу автоматического распознавания образов можно сформулировать на языке теории статистических решений.

Имеются $\pi_1, \pi_2, \dots, \pi_m$ — генеральных совокупностей с плотностями распределения $\rho_1(x), \dots, \rho_m(x)$ и априорными

вероятностями q_1, q_2, \dots, q_m , соответственно. Заданы цены ошибочной классификации $C(\delta_i)$.

Необходимо выбрать решающее правило так, чтобы математическое ожидание потерь

$$P(R, \delta_i) = \sum_{i=1}^m q_i \left\{ \sum_{\substack{j=1 \\ j \neq i}}^m C(\delta_i) \int p_j(x) dx \right\}$$

в фиксированном пространстве $\{R\}$ было минимальным.

Как известно [1], метод Байеса дает минимум математического ожидания потерь в фиксированном $\{R\}$. Он заключается в том, что реализация относится к π_i , если результат измерения x попадает в область R_i . При равных ценах ошибочной классификации $C(\delta_i)$ область R_i состоит из тех точек x , для которых

$$q_i p_i(x) > q_j p_j(x) \quad (j=0,1,\dots,m; j \neq i). \quad (I)$$

Если

$$P \left\{ \frac{p_i(x)}{p_j(x)} = a / \pi_h \right\} = 0,$$

то функция распределения $P_i(x)/P_j(x)$ для любого π_h и совместная функция распределения двух отношений являются непрерывными (вытекает непосредственно из определения образа). Тогда, как показано в работе [1],

- 1) метод Байеса является оптимальным в смысле минимума математического ожидания потерь,
- 2) любой допустимый метод является методом Байеса,
- 3) класс методов Байеса является минимальным полным классом.

2. Получение достаточных оценок $P_i(x)$ по конечной выборке

Метод Байеса предполагает, что распределение всех генеральных совокупностей известны точно. Но при решении большинства практических задач распознавания образов эти распределе-

ния являются неизвестными. В процессе обучения распознавающей системе предъявляется лишь конечная выборка из генеральных совокупностей. Рассмотрим вопрос о том, как можно получить достаточные оценки $P_i(x)$, используя эту выборку.

Ввиду того, что пространство измеряемых системой параметров $\{R\}$ всегда ограничено (из-за ограниченности пределов измерения каждого параметра), а точность измерения их не может быть сколь угодно большой (из-за наличия шумов), в практических задачах распознавания образов пространство $\{R\}$ всегда можно представить в виде конечного числа дискретных точек. Так, если число градаций по x_1, x_2, \dots, x_n равно соответственно $\ell_1, \ell_2, \dots, \ell_n$, то $\{R\}$ содержит в себе $\prod_{i=1}^n \ell_i$ дискретных точек.

Тогда отношение числа реализаций, попавших в точку с координатами x_ξ , к объему выборки N (частота), даст наилучшую (для данной выборки) оценку $P_i(x)$ в этой точке. Проделав эту операцию для всех x , получим гистограмму $\bar{P}_i(x)$, заданную в дискретных точках. Причем, согласно теореме Бернулли, $\bar{P}_i(x)$ стремится к $P_i(x)$ при $N \rightarrow \infty$, т.е.

$$\lim_{N \rightarrow \infty} P(|\frac{n}{N} - P| < \varepsilon) = 1,$$

где $\frac{n}{N}$ — частота события, а P — вероятность (ε — любое положительное число).

Процедуру получения оценок $P_i(x)$ можно осуществить следующим образом: задается пространство $\{R\}$ с координатными осями x_1, x_2, \dots, x_n и с числом градаций по осям $\ell_1, \ell_2, \dots, \ell_n$, соответственно. Каждой точке этого дискретного пространства ставятся в соответствие некоторые ячейки памяти (например, ячейки памяти ЭВМ), число ячеек для представления одной точки пространства равно m (по числу распознаваемых образов). В процессе обучения система предъявляются реализации образов из выборок объемом N_1, N_2, \dots, N_m , и в ячейках памяти суммируется число попавших в эту точку реализаций. Полученная сумма делится на N_i и (при достаточно большом N) получается достаточно хорошая оценка плотностей распределения $P_i(x)$. Таким образом можно в принципе решить любую задачу распознавания образов.

Здесь подчеркивается слово "в принципе", т.к., во-первых, объем выборки при такой процедуре должен быть очень большим и, более того, трудно оценить даже приблизительно его величину; во-вторых, уже при не очень больших n , ℓ и m объем необходимой памяти становится настолько большим, что задача оказывается практически не реализуемой на самых мощных ЭВМ. При $n = 20$, $\ell = 100$, $m = 50$ (размерность многих практических задач распознавания образов), объем необходимой памяти равен $5 \cdot 10^{22}$ ячеек.

Покажем, что этот объем можно существенно сократить, не отходя от принципиальных позиций статистических методов.

Согласно (I), область принятия решения о принадлежности к генеральной совокупности \mathcal{P}_i содержит все те точки x , для которых

$$q_i P_i(x) > q_j P_j(x) \quad (j \neq i).$$

Эта область в общем случае может содержать точки, для которых $q_i P_j(x) = 0$ для всех $j \neq i$ (область уверенных ответов $R_i^{(1)}$). Чтобы принять решения о точках, попавших в эту область, нет необходимости точно знать функцию $P_i(x)$, достаточно лишь знать, что она больше нуля.

С другой стороны, пространство $\{R\}$ может содержать такую область $R_i^{(3)}$, где $P_i(x) = 0$ (хотя бы для некоторых i) и область перекрытий $R_i^{(2)}$, где одновременно $P_i(x) \neq 0$ и $P_j(x) \neq 0$ (хотя бы для одного j). Тогда для каждого i пространство $\{R\}$ можно представить в виде суммы областей:

$$\{R\} = R_i^{(1)} + R_i^{(2)} + R_i^{(3)}.$$

Для принятия решения по (I) достаточно знать функцию $P_i(x)$, заданную лишь на $R_i^{(2)}$ и границы областей $R_i^{(1)}, R_i^{(2)}, R_i^{(3)}$.

Вероятность попадания точки в $R_i^{(2)}$ для реализаций K -ой генеральной совокупности есть

$$q_K \int_{R_i^{(2)}} P_i(x) dx \quad (i=0,1,\dots,m),$$

а вероятность получения от системы в целом неуверенных ответов выражается как

$$\sum_{i=1}^m q_i \int_{R_i^{(2)}} P_i(x) dx \quad (i=0,1,\dots,m), \quad (2)$$

На основании этого определим понятие информативности пространства. Будем говорить, что пространство измеряемых систем признаков $\{R\}$ достаточно информативно для распознавания образов \mathcal{P}_i , если

$$\sum_{i=1}^m q_i \int_{R_i^{(2)}} P_i(x) dx \leq \gamma \quad (i=1,2,\dots,m),$$

где γ - заданная положительная величина в интервале $[0,1]$; и полностью информативно, если

$$\sum_{i=1}^m q_i \int_{R_i^{(2)}} P_i(x) dx = 0.$$

Для получения количественного суждения об информативности пространства $\{R\}$ нужно вычислить (2). В этом случае нет необходимости точно знать $P_i(x)$ в области $R_i^{(2)}$, а достаточно лишь иметь её оценку по конечной выборке объема N_i из генеральной совокупности \mathcal{P}_i , причем величина N_i будет зависеть от задания точности и достоверности вычисления интеграла и не будет зависеть от вида функции распределения вероятности.

Для вычисления интеграла в выражении (2) воспользуемся методом статистических испытаний. Если зафиксирована область $R_i^{(2)}$ (вопрос о получении границ областей $R_i^{(1)}$ и $R_i^{(3)}$ будет рассмотрен ниже) и системе представляются реализации случайного вектора x из выборки объема N , то отношение числа точек, попавших в область $R_i^{(2)}$, ко всему объему выборки даст исключенную оценку [4], причем,

$$\lim_{N \rightarrow \infty} P \left(\left| \frac{N_{R_i^{(2)}}}{N} - \rho \right| < \varepsilon \right) = 1,$$

здесь

$$\rho = \int_{R_i^{(2)}} P_i(x) dx,$$

$n_{R_i}^{(2)}$ - число точек, попавших в область $R_i^{(2)}$,
 ε - сколь угодно малое положительное число,

то есть

$$\frac{n_{R_i}^{(2)}}{N} = \bar{\rho} \approx \rho. \quad (3)$$

Будем говорить, что равенство (3) имеет точность ε с надежностью α , если для неравенства

$$|\rho - \bar{\rho}| < \varepsilon$$

справедливо соотношение

$$P(|\rho - \bar{\rho}| < \varepsilon) = \alpha. \quad (4)$$

Связем точность вычисления ε и надежность α с необходимым объемом выборки N .

Воспользуемся для этого неравенством Чебышева:

$$P(|\bar{\rho} - \rho| < \varepsilon) \geq 1 - \frac{\sigma_{\bar{\rho}}^2}{\varepsilon^2}. \quad (5)$$

Заметим, что $\bar{\rho}$ есть случайная величина с математическим ожиданием

$$M(\bar{\rho}) = \rho$$

и дисперсией

$$D(\bar{\rho}) = \frac{\rho(1-\rho)}{N}.$$

Тогда средняя квадратическая ошибка равенства (3) $\sigma_{\bar{\rho}}$ равна:

$$\sigma_{\bar{\rho}} = \sqrt{\frac{\rho(1-\rho)}{N}}. \quad (6)$$

Нетрудно заметить, что она достигает максимума при $\rho = 0,5$. Сравнивая (4) и (5), запишем:

$$\alpha \geq 1 - \frac{\sigma_{\bar{\rho}}^2}{\varepsilon^2}.$$

Подставляя вместо $\sigma_{\bar{\rho}}^2$ его значение из (6) получим:

$$\alpha \geq 1 - \frac{\rho(1-\rho)}{N\varepsilon^2}$$

или

$$N = \frac{\rho(1-\rho)}{(1-\alpha)\varepsilon^2}, \quad (7)$$

то есть для того, чтобы вычислить интеграл в выражении (2) с точностью ε равной, например 0,01 и надежностью $\alpha = 0,95$ при $\rho = 0,5$, необходимо взять выборку объема $N = 50000$ испытаний.

Формула (7), полученная на основании неравенства Чебышева, дает сильно завышенное значение для N . Более точную оценку для N можно получить в том случае, если использовать закон распределения случайной величины $\bar{\rho}$. Известно [4], что величина $\bar{\rho}$ имеет асимптотически (при $N \rightarrow \infty$) нормальное распределение. На этом основании равенство (4) можно записать в виде:

$$P\left(\frac{|\bar{\rho} - \rho|}{\sigma_{\bar{\rho}}} < t_{\alpha}\right) = \alpha,$$

где t_{α} - величина критического интервала, которая выбирается из таблиц нормального распределения по заданному значению α . Сравнивая это соотношение с (5), получаем:

$$\varepsilon = t_{\alpha} \sigma_{\bar{\rho}}$$

или

$$\varepsilon = t_{\alpha} \sqrt{\frac{\rho(1-\rho)}{N}}.$$

Тогда

$$N = \frac{\rho(1-\rho)}{\varepsilon^2} t_{\alpha}^2. \quad (8)$$

Для $\varepsilon = 0,01$, $\alpha = 0,95$ и $\rho = 0,5$ получим объем необходимой выборки $N = 390$.

3. О представительности выборки

Рассмотрим, как выделяются границы областей $R_i^{(1)}$ и $R_i^{(2)}$ в случае одномерного и многомерного пространства. Вычислим для этого так называемые толерантные (допустимые) пределы.

В одномерном пространстве. Пусть задана желаемая надежность правильной классификации ρ . Это значит, что необходимо выделить область R_i в пространстве $\{R\}$, попадание в которую при единичном испытании не меньше ρ (рис. I) причем достоверность выделения области должна быть не меньше γ , то есть

$$\mathcal{P}\left\{ \int_{R_i} \rho_i(x) dx \geq \rho \right\} \geq \gamma;$$

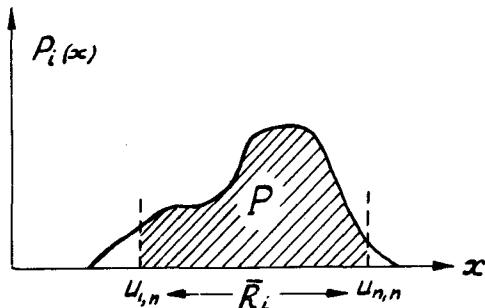


Рис. I. Область \bar{R}_i в одномерном пространстве.

Тогда утверждается [5], что можно подобрать такое N , что крайние члены вариационного ряда, то есть наименьшего $U_{i,n}$ и наибольшего $U_{n,n}$ из выборки объема N , могут быть приняты в качестве границ области \bar{R}_i :

$$U_i(x_1, x_2, \dots, x_n) = U_{i,n},$$

$$U_n(x_1, x_2, \dots, x_n) = U_{n,n}.$$

При этом необходимый объем выборки N определяется из выражения:

$$N\rho^{N-1} - (N-1)\rho^N = 1 - \gamma. \quad (9)$$

Например, для $\gamma = 0,95$ и $\rho = 0,99$ можно найти значение $N \approx 130$. Таким образом, имея выборку объема 130, взятую из произвольной совокупности с любым законом распределения вероятностей, можно утверждать, с достоверностью 0,95, что между крайними членами её лежит не менее 0,99 всей совокупности.

Следовательно, для того, чтобы получить границы области \bar{R}_i , нужно выбрать наибольшее и наименьшее значения из предъявленных реализаций выборки, причем объем выборки зависит от задания количества точности и достоверности и не зависит от вида функции распределения вероятностей образов.

Вычислив границы областей \bar{R}_i , нетрудно получить и гра-

ницы областей $R_i^{(1)}$ и $R_i^{(2)}$

в многомерном пространстве.

для выделения субъекта \bar{R}_i в случае многомерного пространства построим выпуклую оболочку, натянутую на крайние точки из реализаций выборки объема N .

Проведем нелинейное преобразование пространства $x_k \Rightarrow y_k$, так, чтобы выпуклая оболочка области \bar{R}_i отобразилась в κ -мерный параллелепипед \bar{D}_i , грани которого перпендикулярны осям координат (рис.2).

$$y_k = \alpha_k(\varphi) x_k.$$

Наложим дополнительное ограничение – требование неизменности объема области при преобразовании, выражющееся в условии:

$$\prod_{k=1}^{\kappa} \alpha_k(\varphi) = 1.$$

При этом производная от $P_i(x)$ будет иметь конечное число разрывов. Следовательно,

$$\int_{\bar{R}_i} P_i(x) dx = \int_{\bar{D}_i} P_i(y) dy \quad (10)$$

и оба интеграла существуют. Тогда можно сформулировать следующую теорему:

Теорема. Для любых наперед заданных величин ρ и γ , удовлетворяющих условиям $0 < \rho < 1$ и $0 < \gamma < 1$, и независимых параметров y_k , можно указать такой достаточный объем выборки N из генеральной совокупности π , что κ -мерный параллелепипед, натянутый на крайние точки из реализаций этой выборки, ограничит область, вероятность попадания в которую при единичном испытании не меньше ρ с достоверностью не меньше γ , то есть

$$\mathcal{P}\left\{ \int_{\bar{D}_i} P_i(y) dy \geq \rho \right\} \geq \gamma. \quad (II)$$

Доказательство.

Пусть $U_1^{(\kappa)}$ и $U_N^{(\kappa)}$ минимальная и максимальная точки выборки по κ -той оси соответственно. Грани параллелепипеда \bar{D}_i проходят через эти точки перпендикулярно к этим осям.

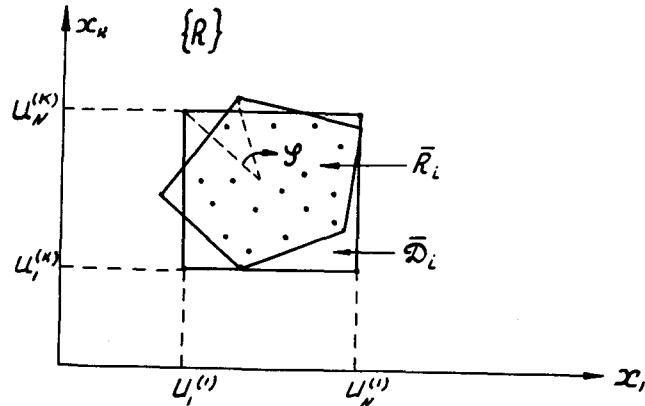


Рис. 2. Проекции областей \bar{R}_i и \bar{D}_i на плоскости $x_k o x_1$.

Тогда

$$\begin{aligned} & \mathcal{P}\left(\int P_i(y) dy > P\right) = \\ & = \mathcal{P}\left\{\prod_{k=1}^N \left[F_i(u_N^{(k)}) - F_i(u_i^{(k)})\right] > P\right\} = \\ & = \mathcal{P}(z_{1,N} \cdot z_{2,N} \cdots z_{n,N} > P), \quad (I2) \end{aligned}$$

где

$$\begin{aligned} F_i(u^{(i)}) &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} P_i(y) dy = \\ & = \mathcal{F}(\infty, \infty, \dots, u^{(k)}, \dots, \infty), \\ z_{k,N} &= F_i(u_N^{(k)}) - F_i(u_i^{(k)}). \end{aligned}$$

Известно, [5] что случайная величина $z = F(y)$ распределена по равномерному закону независимо от дифференциальной функции распределения $P_i(y)$, а разность

$$z_N = F(u_N) - F(u_i)$$

имеет интегральное распределение:

$$\mathcal{P}(z_{1,N} < t) = N \int_0^t [F(u+t) - F(u)]^{N-1} du = N(1-t)^{N-1} + t^N. \quad (0 \leq t \leq 1)$$

Отсюда плотность вероятности случайной величина $z_N^{(k)}$ имеет вид:

$$W(z_{k,N}) = N(N-1)(1-z_{k,N})^{N-2}, \\ z_{k,N} \in [0, 1].$$

Тогда

$$\begin{aligned} \mathcal{P}(z_1 \cdot \dots \cdot z_n > P) &= \int_{z_{1,N}} \int_{z_{2,N}} \cdots \int_{z_{n,N}} W(z_1) \cdot \\ &\cdot W(z_n) dz_1 \cdots dz_n = \\ &= \int_P^\rho W(z_n) dz_n \int_{z_{n-1}}^\rho W(z_{n-1}) dz_{n-1} \cdots \int_{z_2}^\rho W(z_1) dz_1 = \\ &= \psi(N, n, P). \quad (I3) \end{aligned}$$

Приравнивая (I3) к заданной достоверности γ , находим:

$$\psi(N, n, P) = \gamma; \quad (I4)$$

$$N = \varphi(n, P, \gamma).$$

Следовательно, обеспечив объем выборки N из (14), мы получим

$$\mathcal{P}\left\{\left(\frac{1}{N} \sum_{k=1}^N [F_i(\mathcal{U}_N^{(k)}) - F_i(\mathcal{U}_i^{(k)})]\right) > P\right\} \geq r.$$

Тогда в силу (12) выражение (II) справедливо и теорема доказана.

Замечание

В случае зависимых \mathcal{U}_k в большинстве практических задач распознавания величина интеграла в (12) не может быть меньше. Это даёт основание полагать, что теорема будет справедлива и в случае зависимых параметров, хотя строгого доказательства этого утверждения провести не удалось.

На основании теоремы дадим определение представительности выборки. Выборку будем считать представительной, если её объем не меньше N , где N определяется выражением (14).

Итак, имея представительную выборку, можно получить с заданной достоверностью границы областей $R_i^{(1)}$ и $R_i^{(2)}$, внутри которых заключена доля вероятности не меньше P . Это означает, что при дальнейших испытаниях все реализации из генеральной совокупности \mathcal{P}_i будут с вероятностью P попадать внутрь выделенных областей. Если же специфика задачи такова, что объем выборки заранее ограничен, то можно указать достоверность результатов решения задачи распознавания. Задавшись областью перекрытий $R_i^{(2)}$, мы можем вычислить информативность пространства $\{R\}$ (по той же выборке объема N). Продолжать решение задачи имеет смысл тогда, когда информативность пространства нас удовлетворяет; в противном случае следует искать другое, более информативное пространство признаков.

Рассмотрим вопрос, о получении достаточных оценок $P_i(x)$ по конечной выборке в области перекрытий $R_i^{(2)}$. Можно получить гистограмму $\bar{P}_i(x)$ в этой области вычислением величины $\frac{N_x}{N}$, где N_x — число точек, попавших в элементарный гипперобъем; но для этого нужно иметь большой объем памяти и большой объем выборки, причем вопрос о достаточном объеме выборки опять остается нерешенным.

Как уже говорилось, вероятность попадания реализации из генеральной совокупности \mathcal{P}_i в область перекрытия $R_i^{(2)}$ выражается выражением:

$$q_i \int_{R_i^{(2)}} P_i(x) dx.$$

Величину интеграла можно вычислить с заданной точностью и достоверностью, имея конечную выборку объема N . При условии достаточной информативности пространства области перекрытий, как правило, будут содержать "хвосты" функций распределения плотностей вероятности образов. Эти "хвосты" независимо от вида всей функции, могут быть достаточно хорошо аппроксимированы функцией κx^t . Тогда, положив $P_i(x)$ на границе области $R_i^{(1)}$ равной нулю и зная величину интеграла в (2), получим значение K из уравнения:

$$\kappa \int_{R_i^{(2)}} x^t dx = c,$$

где c — значение интеграла в (2), вычисленного методом статистических испытаний. При этом величина средней квадратической ошибки отклонения $\bar{P}_i(x)$ от $P_i(x)$ будет не больше ε . Тогда значение $\bar{P}_i(x)$ в области $R_i^{(2)}$ получим из уравнения

$$\bar{P}_i(x) = \kappa x^t \quad (16)$$

(величина t может варьироваться).

Таким образом, в процессе обучения мы получили оценку информативности выбранного пространства параметров $\{R\}$, границы областей уверенных ответов $R_i^{(1)}$, границы областей перекрытий $R_i^{(2)}$ и приближение с заданной точностью $\bar{P}_i(x)$ в области $R_i^{(2)}$ к функции плотности распределения вероятностей образов $P_i(x)$.

В процессе распознавания в зависимости от специфики задачи могут быть выбраны различные стратегии: стратегия Байеса, метод последовательного анализа Вальда, стратегия неуверенных ответов. Рассмотрим подробнее каждую из них.

Стратегия Байеса. Предъявляя системе неизвестную реализацию для распознавания (т.е. задавая значения x_ξ в пространстве $\{R\}$, определяем, в какую из возможных областей $R_i^{(1)}$ или $R_i^{(2)}$ она попадает. Если точка попадет в область $R_i^{(1)}$, то принимается решение о принадлежности предъявленной реализации к генеральной совокупности \mathcal{P}_i , если же она по-

падет в область $R_i^{(2)}$, то по (16) вычисляется значение $\bar{P}_i(x)$ в этой точке для всех i и выбирается тот индекс i , для которого $q_i \bar{P}_i(x)$ максимально. При попадании точки в область $R_i^{(3)}$ (область, где $P_i(x) = 0$) принимается гипотеза π_0 . При этом математическое ожидание потерь будет:

$$\sum_{i=1}^m q_i \sum_{\substack{j=1 \\ j \neq i}}^m \int_{R_j^{(2)}} P_i(x) dx + \sum_{i=1}^m (1 - P_i), \quad (17)$$

где P_i — доля вероятности, заключенная в области \bar{R}_i . Интеграл в (17) вычисляется методом статистических испытаний в процессе обучения по конечной выборке объема N .

Для сокращения времени при распознавании определение принадлежности точки к той или иной области делается последовательно за n шагов методом дихотомического деления пространства [9].

Метод последовательного анализа Вальда. Системе предъявляется неизвестная реализация, определяется, какой из возможных областей $R_i^{(1)}$ или $R_i^{(2)}$ принадлежит точка. Если точка попадает в область $R_i^{(1)}$, то принимается гипотеза о принадлежности реализации к генеральной совокупности π_i , если же точка попадает в область неуверенных ответов, то вычисляется отношение правдоподобия

$$\lambda_{kij} = \frac{\bar{P}_{ik}(x)}{\bar{P}_{jk}(x)} \quad (18)$$

для тех индексов i и j , для которых $\bar{P}_i(j)(x)$ в найденной области отличны от нуля.

Если $\lambda_{kij} \geq A$, то принимается гипотеза о принадлежности реализации к π_i ; если $\lambda_{kij} \leq B$, то реализация относится к π_j , и если $B_{ij} < \lambda_{kij} < A_{ij}$, то принимается решение о повторном испытании (переспросе) для вычисления следующего значения $\lambda_{(k+1)ij}$. Переспрос продолжается до тех пор, пока не будет принята гипотеза о принадлежности x к тому или иному образу. Константы A_{ij} и B_{ij} ($A > B$) связаны с вероятностями ошибок первого и второго рода α и β следующим соотношением:

$$A \leq \frac{1-\beta}{\alpha}; \quad (19)$$

$$B \geq \frac{\beta}{1-\alpha}.$$

Для независимых и имеющих одну и ту же функцию распределения наблюдений \bar{P}_{ik} определяется как:

$$\bar{P}_{ik} = \prod_{n=1}^K \bar{P}_i(x_n) \quad (n=1, 2, \dots, K).$$

Значения $\bar{P}_i(x_n)$ вычисляются из (16).

Вальд и Вольфович [7] показали, что метод последовательного анализа оптимальен в смысле минимизации среднего риска, получаемого в случае единичного испытания; кроме того, для заданных α и β , он минимизирует среднее число переспросов. Следует заметить, что в том случае, когда на каждом K -ом шаге точка x_K будет попадать в область равных вероятностей конкурирующих гипотез, метод Вальда не будет сходиться. Поэтому при решении практических задач распознавания образов имеет смысл применить так называемый усеченный метод последовательных испытаний. Усечение заключается в том, что заранее назначается максимально допустимое число переспросов n . Последовательные испытания проводятся до тех пор, пока не будет принято решение или не будет проверено n -ое наблюдение. Если при n -ом наблюдении решение еще не принято, то переходят либо к стратегии Байеса, либо к стратегии неуверенных ответов, принимая за измерение реализации значение λ_n .

Стратегия неуверенных ответов в отличие от вышеописанных стратегий, при попадании точки x в область неуверенных ответов вычисляются значения $\bar{P}_i(x)$ по (16), и система выдает те индексы i , со значениями их вероятностей в этой точке, для которых $\bar{P}_i(x) > 0$, [8].

Основные выводы

1) Задача распознавания образов по своей сущности является статистической и может быть решена статистическими методами независимо от того, априори известны или нет функции распределения образов. Достоверность результатов решения задачи зависит от объема выборки обучающей последовательности.

2) Задача распознавания решается в несколько этапов: задаются алфавит распознаваемых образов, желаемая надежность и достоверность решения задачи, выбирается ориентировочно пространство параметров и оценивается информативность этого пространства в процессе обучения; если информативность пространства не удовлетворяет поставленным требованиям, то ведется поиск но-

вого пространства до тех пор, пока не будут удовлетворены эти требования (вопрос о выборе пространства параметров является самостоятельным и здесь не рассматривается).

3) Алгоритм решения задачи заключается в следующем: задаются: алфавит распознаваемых образов $I, 2, \dots, m$; пространство измеряемых параметров $I, 2, \dots, n$; число градаций по каждому из параметров $\ell_1, \ell_2, \dots, \ell_n$, доля вероятности P_i , заключенная в области R_i , и достоверность выделения границы области Γ . Вычисляется по (14) необходимый объем выборки N_i . В процессе обучения для всех m образов строятся границы областей R_i . Определяются границы областей $R_i^{(1)}$ и $R_i^{(2)}$ и вычисляется информативность пространства по (2). В зависимости от специфики задачи выбирается одна из возможных стратегий: стратегия Бейеса, метод последовательного анализа Вальда или стратегия неуверенных ответов. Если выбрана стратегия Вальда, то дополнительно задаются значения α_{ij} и β_{ij} и тогда по (19) вычисляются значения A_{ij} и B_{ij} . Если выбрана стратегия Бейеса, то вычисляется по (17) математическое ожидание потерь.

После этого процесс обучения считается временно законченным и система готова к непосредственному распознаванию образов. При предъявлении системе неизвестной реализации определяется, к какой из областей $R_i^{(1)}$ или $R_i^{(2)}$ она относится, и выдается i -й индекс, если реализация попадает в область $R_i^{(1)}$. Если реализация попадает в область $R_i^{(2)}$, то либо вычисляется i -й индекс по (1) (в зависимости от выбранной стратегии), либо принимается решение согласно (18), либо выдаются все конкурирующие гипотезы с их значениями вероятностей в этой области.

4) В процессе работы системы возможно её совершенствование (подучивание), при этом изменяются границы областей $R_i^{(1)}$ и $R_i^{(2)}$ и повышается достоверность результатов распознавания.

5) Задачу распознавания образов в общем случае нужно рассматривать как "многоэтажную", иерархическую систему. На первом этапе распознаются образы в пространстве физических, химических или каких-либо других свойств предметов и явлений; на следующем этапе распознаются образы в пространстве признаков, которые состоят из алфавита распознанных образов и т.д.

Примером такой иерархической системы может служить задача распознавания речи. На первом этапе проводится членение непрерывного речевого потока на фонемы; затем распознаются фонемы на выделенных участках, далее в пространстве фонем распознают-

ся слова и т.д.

В статье рассмотрена элементарная система для решения задачи распознавания образов, элементарная в том смысле, что она может служить элементом построения иерархической системы любой сложности.

Автор считает своим приятным долгом выразить глубокую благодарность Н.Г. Загоруйко, Л.Я. Савельеву, Э.Х. Гимадутдинову, Г.Я. Волошину и Г.С. Лбову, оказавшим помочь в настоящей работе.

ЛИТЕРАТУРА

1. Т. Андерсон. Введение в многомерный статистический анализ. М., ФМ, 1963, стр. 175-211.
2. К. Фу. Модель последовательных решений для оптимального опознавания. - Проблемы бионики. М., "Мир", 1965., стр. 374-384.
3. В.А. Ковалевский. Задача распознавания образов с точки зрения математической статистики. - Семинар "Распознавание образов и конструирование читающих автоматов". Киев, 1965.
4. Н.П. Бусленко, Ю.А. Шрейдер. Метод статистических испытаний. М., ФМ, 1961, стр. 66-72.
5. И.В. Дунин-Барковский, Н.В. Смирнов. Теория вероятностей и математическая статистика в технике. М., Гостехиздат, 1955.
6. В.Н. Елкина, Н.Г. Загоруйко. Алгоритм поиска формальных элементов алфавита. (Данный сборник, стр.59).
7. Wald A., Wolfowitz J., Optimum character of the sequential probability ratio test. Ann. Math. Stat., 19, 326 - 339, 1948.
8. Streck G.P. Stochastic Model for the Browning-Bledsoe Pattern Recognition Scheme. IRE Trans. on Electr. Comp. 1962, EC-11, pp.274-282.
9. Forgie J.W. and Forgie C.D. Computer Identification of Vowel Types. GASA, vol.33, 1, January 1961, p.7-11.

Поступила в редакцию
11.1.1966.