

ВЫЧИСЛИТЕЛЬНЫЕ СИСТЕМЫ

Сборник трудов
Института математики СО АН СССР

1966 г.

Выпуск 22

О ПРЕДСТАВЛЕННОСТИ ВЫБОРКИ ПРИ ВЫБОРЕ
ЭФФЕКТИВНОЙ СИСТЕМЫ ПРИЗНАКОВ

Г.С. Ибов

§ I. Постановка задачи

Как показано в работе [1], эффект, вызываемый ограниченностью экспериментального материала, может оказаться весьма существенное влияние на точность определения вероятности неправильного распознавания, что, в свою очередь, значительно усложняет выбор эффективной системы признаков [2].

В настоящей работе даны некоторые рекомендации по выбору эффективной системы признаков в условиях ограниченной выборки. В работе также приводятся результаты эксперимента на ЭВМ по установлению достаточного объема выборки для определения вероятности неправильной классификации с заданной точностью (случай нормального распределения).

Рассмотрим один из возможных подходов к решению задачи выбора эффективной системы признаков в случае неизвестных законов распределения. Делаем произвольное предположение о типе распределения генеральных совокупностей $\rho'_1(x), \dots, \rho'_K(x)$, соответствующих K образам. Здесь через x обозначен вектор в пространстве признаков, на которых строится рассматриваемое про-

пространство. Разбиваем пространство признаков на непересекающиеся области U'_1, \dots, U'_K (по числу образов). Область U'_i определяем как совокупность тех точек x , для которых

$$q_i P'_i(x) = \sup_j \{q_j P'_j(x)\}_{j=1,\dots,K}. \quad (1)$$

Положим, что $P'_1(x), \dots, P'_K(x)$ — нормальные распределения. Это означает, что мы будем использовать решающие функции либо из класса линейных, либо квадратичных функций. Выборку, на основе которой оцениваем параметры распределений $P'_1(x), \dots, P'_K(x)$, назовем обучающей выборкой. При этом, очевидно, решающее правило (1) будет отличаться от оптимального решающего правила, во-первых, из-за возможного неправильного предположения о типе закона распределения и, во-вторых, из-за использования оценок параметров распределения по ограниченной обучающей выборке вместо самих параметров распределения.

Обозначим через φ' вероятность неправильного распознавания, получаемую при условии оптимального разбиения пространства [6].

В качестве оценки вероятности неправильной классификации можем использовать либо функционал φ^* [2], либо, если существует возможность получить контрольную выборку, величину φ' , определяемую следующим образом:

$$\varphi' = \sum_{i=1}^K q_i \sum_{\substack{j=1 \\ j \neq i}}^K P'(j/i) = \sum_{i=1}^K q_i \varphi'_i, \quad (2)$$

где

$$P'(j/i) = \int_{U'_i} P_j(x) dx. \quad (3)$$

Заметим, что под интегралом используется истинное распределение $P_j(x)$, а не $P'_j(x)$. Другими словами, значение интеграла (3) оценивается по контрольной выборке после того, как определены области U'_1, \dots, U'_K согласно решающему правилу (1). Априорные вероятности q_1, \dots, q_K будем считать известными. Замечания о необходимом объеме контрольной выборки будут приведены ниже.

Обозначим через N_i число реализаций, составляющих обучающую выборку для i -го образа. Если увеличивать объем

обучающей выборки ($N_1 \rightarrow \infty, \dots, N_K \rightarrow \infty$), то величина φ' будет стремиться по вероятности к некоторому значению P_1 , вообще говоря, отличному от φ . Равенство величин P_1 и φ достигается лишь при условии правильного предположения о нормальном законе распределения. Величина отклонения P_1 от φ зависит от того, насколько хорошо неизвестные распределения $P_1(x), \dots, P_K(x)$ аппроксимируются распределениями $P'_1(x), \dots, P'_K(x)$. Решающие функции, при которых вероятность неправильной классификации равна P_1 , будут наилучшими, дающими минимальную вероятность ошибки для выбранного класса решающих функций. Заметим, что при ограниченном объеме обучающей выборки будет всегда выполняться соотношение $\varphi' \geq P_1$. Величина φ' является случайной величиной, так как представляет собой функцию случайных наблюдений, составляющих обучающую выборку. Обозначим плотность распределения φ' через $\varphi(\varphi')$.

Поясним изложенное с помощью рис. I.

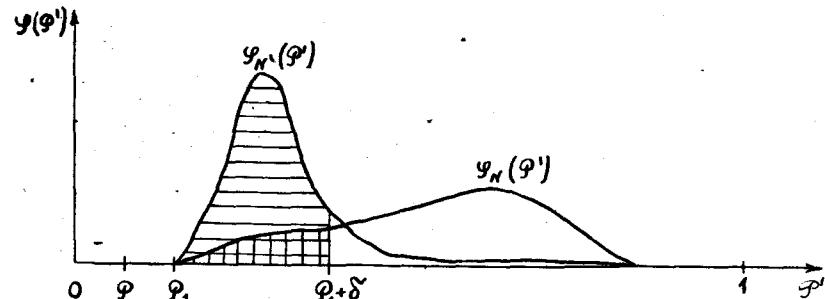


Рис. I. Плотности распределения величины φ' для двух различных объемов обучающей выборки.

Пусть $\varphi_N(\varphi')$ отображает плотность распределения φ' при некотором фиксированном объеме выборки $N_1 = N_2 = \dots = N_K = N$. Функция $\varphi_{N'}(\varphi')$ отображает плотность распределения φ' при объеме выборки $N'_1 = N'_2 = \dots = N'_K = N'$, причем $N' > N$. Здесь δ — погрешность в определении P_1 , вызванная ограниченностью обучающей выборки.

При увеличении объема обучающей выборки величина φ' стремится к P_1 и увеличивается вероятность получения неравенства $\varphi' < P_1 + \delta$ при фиксированном значении δ (см. заштрихованные области на рис. I для плотностей $\varphi_N(\varphi')$ и $\varphi_{N'}(\varphi')$). Кроме того, оказывается, что при фиксированном объеме обучающей выборки и при данном значении P_1 вероятность равного отклонения φ'

от P_1 тем больше, чем больше размерность пространства (§8). На практике часто встречается случай, когда при высокой размерности исходной системы признаков имеем ограниченный по тем или иным причинам объём выборки. Это может привести к существенному отклонению величины \bar{P}' от P_1 . В этих условиях мы можем либо "забраковать" систему признаков, в то время как для этой системы $P_1 < P_0$ (P_0 - допустимая вероятность неправильной классификации), либо при сравнении двух систем X и Y получить $\bar{P}'_X > \bar{P}'_Y$, тогда как на самом деле $P_X < P_Y$. Решение задачи выбора эффективной системы признаков в этих условиях значительно усложняется. Для её решения потребуется понятие представительности выборки.

Для фиксированной системы признаков X и данного комплекса условий D получения реализаций признаков назовём обучавшую выборку представительной, если выполняется следующее неравенство:

$$P_2 \{(\bar{P}' - P_1) < \delta_0\} \geq \beta, \quad (4)$$

где δ_0 - допустимая погрешность в определении вероятности ошибки P_1 (значение β близко к единице).

Величину \bar{P}' (см. выражение (2)) определяем по контрольной выборке. Объём контрольной выборки [3] для определения \bar{P}' с погрешностью ε задаётся выражением:

$$N_o = \frac{\sum_{i=1}^k q_i P'_i (1-P'_i)}{\varepsilon^2} t_\alpha^2, \quad (5)$$

где t_α - величина критического интервала, которая выбирается из таблиц нормального распределения по заданной достоверности α . При этом число контрольных реализаций N_{oi} , принадлежащих i -му образу, равно $q_i N_o$. Получаем

$$P_2 \{|\bar{P}' - P'| < \varepsilon\} \geq \alpha, \quad (6)$$

где \bar{P}' - оценка P' по ограниченной контрольной выборке объёма N_o . Так как вероятности P'_1, \dots, P'_K неизвестны, то для определения N_o поступают следующим образом: задаются либо максимальным значением N_o , которое получается при $P'_1 = P'_2 = \dots = P'_K = 0,5$, либо задаются некоторым первоначальным значением N_o , а затем, найдя оценки $\bar{P}'_1, \dots, \bar{P}'_K$, уточняют значение

числа контрольных реализаций. При этом

$$\bar{P}'_i = \frac{N_{oi}}{N_{oc}},$$

где \bar{N}_{oc} - число неправильно распознанных реализаций i -го образа.

В следующем параграфе рассмотрим решение задачи выбора эффективной системы признаков в условиях ограниченной обучавшей выборки.

§ 2. Выбор эффективной системы признаков.

Пусть задана некоторая исходная система признаков x_1, \dots, x_n . Строим на указанных признаках подпространства n -мерной размерности. При каждом m ($m=1, 2, \dots, n$) будем рассматривать только самое информативное подпространство ^{*)} из всех m -мерных подпространств, число которых равно C_n^m . Число всех различных подпространств равно $2^n - 1$.

Допустим, что при каждом значении m для самого информативного подпространства известны вероятности P_m и $P_m + \delta$, причем величина δ выбрана такой, что

$$P_2 \{(\bar{P}' - P_m) < \delta\} \geq \beta. \quad (7)$$

Обозначим $P_m + \delta$ через \bar{P}_m . На рис. 2 функции $P_m(m)$ и $\bar{P}_m(m)$ заданы в точках, но для наглядности их значения для разных m соединены сплошными линиями.

Через $P(0)$ обозначена вероятность ошибки, получаемая при использовании лишь априорных сведений об образах. Величина $P(0)$ равна $1 - q_i$, если $q_i > q_j$ ($i=1, \dots, K; j \neq i$).

Функции $P_m(m)$ и $\bar{P}_m(m)$ для каждого m задают интервал, внутри которого с большой вероятностью β окажется выборочное значение $P_m(m)$ при данном объёме обучавшей выборки. Такой интервал указан только для самого информативного подпространства при каждом значении m . Для других подпространств n -мерной размерности величина этого интервала тем больше, чем большее вероятность ошибочной классификации P (§8). Если, однако, значение вероятности P приближается к $P(0)$, то величина интервала δ с некоторого момента будет меньше допустимой.

^{*)} Из двух подпространств будем считать более информативным то подпространство, использование которого приводит к меньшей вероятности ошибки.

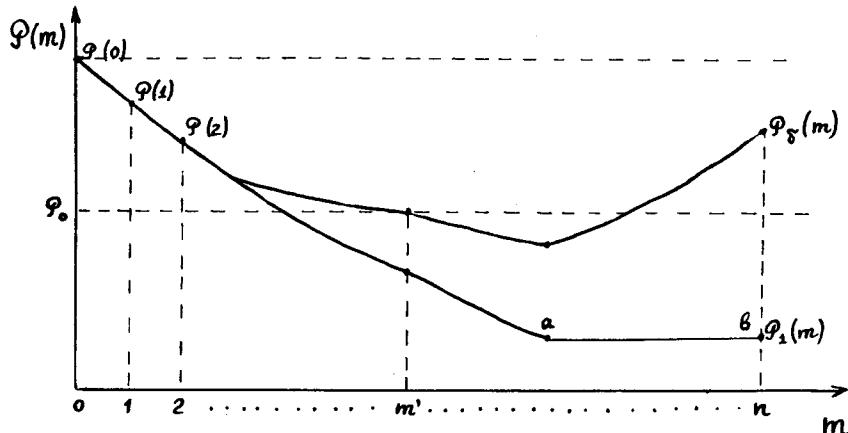


Рис. 2. Доверительный интервал величины \mathcal{P}' для самых информативных подпространств m -мерной размерности.

мого значения δ_0 . Для такого подпространства в соответствии с (4) выборка любого объёма будет представительной. Это объясняется тем, что при $\mathcal{P} = \mathcal{P}(0)$ (пространство используемых признаков не несёт информации) вероятность ошибки не меняется от способа разбиения пространства на области U'_1, \dots, U'_K .

Если функция $P_1(m)$ остается постоянной при увеличении m (область ab рис.2), то функция $P_0(m)$ будет возрастающей. Это может привести к тому, что в некоторых конкретных задачах использование всей системы признаков при фиксированном объёме обучающей выборки может привести к вероятности неправильной классификации большей, чем вероятность неправильной классификации в случае использования лишь части информативных признаков этой системы. При данном объёме выборки N найдется такое m' , что для всех подпространств m -мерной размерности ($m < m'$) величина интервала δ не превзойдёт некоторой наперёд заданной δ_0 . Тогда как для $m > m'$ величина интервала δ может быть больше δ_0 . В этом случае выбор эффективного m -мерного подпространства для $m > m'$ из множества m -мерных подпространств может быть неудачным: для $m > m'$ невозможно установить, насколько предпочтительнее одно подпространство перед другим, так как разница между информа-

тивностями, получаемыми при использовании выбранных двух подпространств, может оказаться меньше, чем погрешность при определении самих информативностей этих подпространств. В описанных условиях для решения задачи выбора эффективной системы признаков можно предложить следующий подход. Сначала выбираем наиболее информативное подпространство m' -мерной размерности. Если окажется, что для него $\mathcal{P}'(m') < \mathcal{P}_0$, то задача может считаться решенной. Если $\mathcal{P}'(m') > \mathcal{P}_0$, то можно увеличивать число признаков исходной системы каждый раз на некоторую величину, оставляя постоянным число реализаций (N_1, \dots, N_K) . При каждом таком увеличении выбирать наиболее информативное подпространство m' -мерной размерности. Стремление к увеличению размерности исходного пространства может быть оправдано тем, что при увеличении n увеличивается общее число m -мерных подпространств, из которых выбирается самое информативное подпространство. Это приводит к увеличению вероятности получения такого подпространства, при котором достигается решение задачи. Однако может получиться так, что все известные признаки исчерпаны, но не оказалось такого m' -мерного подпространства, при котором получаем решение задачи, то есть $\mathcal{P}'(m') < \mathcal{P}_0$. Поэтому, по-видимому, лучше всего, увеличивая n , одновременно увеличивать число реализаций (N_1, \dots, N_K) . Задача может быть не решена только в случае, когда исчерпаны все известные признаки и m' стало равным n , но при этом $\mathcal{P}(n) > \mathcal{P}_0$.

Можно, конечно, с самого начала составить исходную систему из всех известных признаков и обеспечить представительность обучающей выборки для выбранной системы признаков. Однако такой подход менее приемлем, поскольку решение задачи при вышеописанном подходе может быть достигнуто раньше, чем выборка станет представительной для исходной системы из всех известных признаков.

При предлагаемом подходе необходимо использовать алгоритм сокращения числа признаков с n до m' , минуя промежуточные значения m ($m > m'$). В описанных условиях алгоритм, предлагаемый в работе [4], применять нежелательно, так как этот алгоритм предполагает на первом этапе выбор наиболее информативного подпространства среди $(n-1)$ -мерных подпространств, на втором этапе - среди $(n-2)$ -мерных подпространств и т. д. В этом случае можно использовать алгоритм случайного поиска с адаптацией, приведенной в работе [2].

§ 3. Представительность выборки из нормального распределения (случай дихотомии)

В настоящее время при построении решающего правила по выборке руководствуются следующим принципом: чем больше объем обучающей выборки, тем лучше. Такой подход дает лишь интуитивную уверенность в том, что объем используемой выборки достаточночен для оценки вероятности неправильной классификации.

Задача об установлении представительности выборки заданного объема достаточно сложная.

В работе [5] в качестве меры представительности выборки введена величина $\Delta = M\varphi' - \bar{\varphi}$ ($M\varphi'$ – математическое ожидание случайной величины φ') и установлена зависимость величины Δ от объема выборки для одномерного нормального распределения с известными параметрами. Попытка выяснить эту зависимость для многомерного распределения приводит к интегралу, который не может быть выраженным через элементарные функции.

В случае неизвестного распределения можно предложить следующий способ определения представительности выборки заданного объема N ($N_1 = N_2 = \dots = N_k = N$). Получаем выборку объема, достаточного для того, чтобы образовать из нее γ различных групп по N реализаций в каждой группе.

Так как значение φ_1 неизвестно, мы не можем воспользоваться выражением (4) для определения представительности выборки. В этом случае поступим следующим образом. Каждую группу реализаций рассматриваем как отдельную обучающую выборку. С помощью выражения (2) вычисляем вероятность φ' для каждой группы ($\varphi'_1, \dots, \varphi'_k$) и определяем из $\varphi'_1, \dots, \varphi'_k$ минимальное φ'_{min} и максимальное φ'_{max} значения, которые образуют некоторый интервал ($\varphi'_{min}, \varphi'_{max}$). Число γ должно быть таким, чтобы

$$\Pr\{\varphi' \in (\varphi'_{min}, \varphi'_{max})\} > \beta. \quad (8)$$

В данном случае в качестве меры представительности выборки можно выбрать величину $\Delta\varphi' = \varphi'_{max} - \varphi'_{min}$. Если при этом величина интервала $\Delta\varphi'$ меньше некоторой допустимой величины $\Delta\varphi'_0$, то выборку будем считать представительной. Выбор $\Delta\varphi'$ в качестве меры представительности выборки может быть оправдан тем, что при увеличении объема обучающей выборки N выборочные значения φ' все теснее группируются около φ_1 и величина интервала $\Delta\varphi'$ по вероятности стремится к нулю при фиксированном значении γ .

Был поставлен следующий эксперимент для определения представительности выборки объема N из многомерных нормальных распределений. С помощью датчика [3] вырабатывались случайные векторы из двух генеральных совокупностей π_1 и π_2 с некоторыми плотностями распределения $P_1(x), P_2(x)$, приближенно описываемыми нормальными законами распределения $N_1(\mu^{(1)}, E)$ и $N_2(\mu^{(2)}, E)$, где

$$\mu^{(1)} = \begin{Bmatrix} 0 \\ \vdots \\ 0 \end{Bmatrix} \quad \text{и} \quad \mu^{(2)} = \begin{Bmatrix} \frac{\rho}{\sqrt{n}} \\ \vdots \\ \frac{\rho}{\sqrt{n}} \end{Bmatrix}.$$

Величина ρ – евклидово расстояние между центрами нормальных распределений; E – единичная матрица ковариации. Диаграмма для нормальных распределений с указанными параметрами

$$(\mu^{(1)} - \mu^{(2)})' E (\mu^{(1)} - \mu^{(2)}) = \rho^2$$

и вероятность неправильной классификации [6]

$$\bar{\varphi} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}y^2} dy. \quad (9)$$

Для выяснения зависимости величины интервала $\Delta\varphi'$ от информативности и размерности пространства, а также объема обучающей выборки необходимо эту величину определить для некоторых значений $\bar{\varphi}, n, N$ (в случае многомерных нормальных распределений с указанными параметрами). Величина интервала $\Delta\varphi'$ была вычислена при следующих значениях $\bar{\varphi}$: $\bar{\varphi}_1 = 0,0099 \approx 0,01$ ($\rho_1 = 4,66$, см. таблицу I приложения), $\bar{\varphi}_2 = 0,1003 \approx 0,1$ ($\rho_2 = 2,56$, см. таблицу 2 приложения), $\bar{\varphi}_3 = 0,2981 \approx 0,3$ ($\rho_3 = 1,06$, см. таблицу 3 приложения), $\bar{\varphi}_4 = 0,4483 \approx 0,45$ ($\rho_4 = 0,13$, см. таблицу 4 приложения). При каждом значении $\bar{\varphi}$ величина $\Delta\varphi'$ определялась для двух-, шести- и десятимерного пространства и для объема выборки $N = 20, 50, 100, 150$. Таким образом, значение $\Delta\varphi'$ было определено в 48 случаях.

В каждом таком случае находились выборочные значения $\varphi'_1, \dots, \varphi'_N$

случайной величины $\bar{\mathcal{P}}'$ ($\tau=15$). При этом, например, величина $\bar{\mathcal{P}}'_j$ ($j=1, \dots, 15$) определялась так: для данной размерности n пространства признаков случайным образом получали N реализаций из генеральной совокупности π_1 и N реализаций из генеральной совокупности π_2 . Определялись оценки $\bar{\mu}^{(1)}$ и $\bar{\mu}^{(2)}$

$$\bar{\mu}^{(1)} = \begin{Bmatrix} \bar{x}_1^{(1)} \\ \vdots \\ \bar{x}_n^{(1)} \end{Bmatrix}, \quad \bar{\mu}^{(2)} = \begin{Bmatrix} \bar{x}_1^{(2)} \\ \vdots \\ \bar{x}_n^{(2)} \end{Bmatrix}.$$

После этого считалось, что обучение закончено и в соответствии с (1) устанавливалось следующее решающее правило: область Y_i ($i=1, 2$) определялась как совокупность тех точек x , для которых

$$-2 \sum_{p=1}^n \bar{x}_p^{(i)} x_p + \sum_{p=1}^n (\bar{x}_p^{(i)})^2 = \inf_{\ell} \left\{ -2 \sum_{p=1}^n \bar{x}_p^{(\ell)} x_p + \sum_{p=1}^n (\bar{x}_p^{(\ell)})^2 \right\}_{\ell=1, 2}.$$

Для определения $\bar{\mathcal{P}}'$ использовалась контрольная выборка объема $N_{01} = 3000$ из генеральной совокупности π_1 и $N_{02} = 3000$ из π_2 . Если положить $q_1 = q_2 = \frac{1}{2}$, $t_\alpha = 1,96$ и $\bar{\mathcal{P}}'_j \approx \bar{\mathcal{P}}$, то точность вычисления $\bar{\mathcal{P}}'$ согласно (5) будет:

$$\varepsilon = 1,96 \sqrt{\frac{\bar{\mathcal{P}}(1-\bar{\mathcal{P}})}{2 \cdot 3000}}.$$

Значения ε и $\bar{\mathcal{P}}$ приводятся в нижеследующей таблице:

$\bar{\mathcal{P}}$	0,01	0,1	0,3	0,45
ε	0,0002	0,0074	0,0110	0,0130

Число τ выборочных значений величины $\bar{\mathcal{P}}'$ задавалось следующим образом. Известно [7], что для произвольного закона распределения достоверность γ получения неравенства (8) и значения β и τ связаны связью соотношением

$$\gamma = 1 - \tau \beta^{\tau-1} + (\tau-1) \beta^\tau.$$

Например, для $\tau = 15$ и $\gamma = 0,98$ получаем $\beta = 0,9$.

На основе таблиц I, 2, 3, 4 приложения составлена таблица. Если для данных $\bar{\mathcal{P}}$, n , N величина полученного интервала $\Delta \bar{\mathcal{P}}'$ превышала некоторую заданную величину $\Delta \bar{\mathcal{P}}_0$, то считалось, что выборка объема N непредставительна. Этим значением $\bar{\mathcal{P}}, n, N$ в таблице соответствуют нули. Если при некоторых значениях $\bar{\mathcal{P}}$, n , N величина $\Delta \bar{\mathcal{P}}'$ меньше $\Delta \bar{\mathcal{P}}_0$, то этим значениям в таблице соответствуют единицы (выборка объема N считалась представительной).

Из таблицы видно, что выборка одного и того же объема становится более представительной при увеличении информативности пространства. Выборка становится менее представительной с увеличением размерности пространства при фиксированных значениях N и $\bar{\mathcal{P}}$. Те случаи, которые не соответствуют этому правилу, обозначены звездочкой. Это несоответствие, по-видимому, объясняется ограниченностью числа τ выборочных значений случайной величины $\bar{\mathcal{P}}'$.

Полученные результаты можно использовать, когда априори известно, что генеральные совокупности π_1 и π_2 имеют распределения, хорошо аппроксимируемыми нормальными распределениями с единичными матрицами ковариации.

На основе полученных данных можно дать некоторые рекомендации. Если, например, дана допустимая вероятность неправильной классификации $\bar{\mathcal{P}}_0$ и размерность n пространства, то можно оценить приближенный объем выборки, необходимый для того, чтобы не "забраковать" систему признаков X при $\bar{\mathcal{P}}_X < \bar{\mathcal{P}}_0$.

Пусть, например, $\bar{\mathcal{P}}_0 = 0,12$ и $n = 10$. Если при этом $\bar{\mathcal{P}}_X = 0,10$, то необходимо использовать выборку $N = 50$. Действительно, для $N = 50$ $\text{Вер}\{\bar{\mathcal{P}}'_X > \bar{\mathcal{P}}_0\}_{N=50} \approx 0$, в то время как для $N = 20$ $\text{Вер}\{\bar{\mathcal{P}}'_X > \bar{\mathcal{P}}_0\}_{N=20} > 0$.

Далее, если для двух m -мерных подпространств X и Y ($m \leq 10$) будет соблюдаться неравенство $\bar{\mathcal{P}}_Y - \bar{\mathcal{P}}_X > 0,04$, то при $N = 20$ $\text{Вер}\{\bar{\mathcal{P}}'_Y > \bar{\mathcal{P}}'_X\}_{N=20} \approx 0$.

В заключение сделаем ряд следующих замечаний.

I. Величина $\bar{\mathcal{P}}$, определяемая выражением (9), не равна вероятности неправильной классификации $\bar{\mathcal{P}}_1$, поскольку истинные математические ожидания распределений случайных чисел, вырабатываемых датчиком, вообще говоря, отличны от предполагаемых значений $\bar{\mu}^{(1)}$ и $\bar{\mu}^{(2)}$. По этой причине выборочные значения величин-

ны \mathcal{P}' (см. таблицы I, 2, 3, 4 из приложения) могут получиться иногда меньше вероятности \mathcal{P} .

2. Если предположить, что $\varphi'_{min} - \varphi_1 \ll \Delta\varphi'$, то интервал $\Delta\varphi'$ можно считать равным интервалу δ .

3. Если вместо вероятности \mathcal{P}' , определяемой по формуле (2), использовать функционал $\mathcal{P}^*[2]$, то значение \mathcal{P}^* , как правило, будет меньше вероятности неправильной классификации (в то время, как величина \mathcal{P}' всегда больше или равна \mathcal{P}_2). Желательно было бы определить величину интервала δ для \mathcal{P}^* .

4. Если имеем обучаемую выборку объёма N и контрольную выборку объёма N_0 , то $\text{Вер}\{\bar{P} > P_1 + \delta + \varepsilon\} \approx 0$ при значениях α и β , близких к единице (см. выражения (6) и (7)).

5. Проведенный эксперимент можно продолжить для произвольных значений \mathcal{P} , α , N и числа образов $K > 2$. Причем при $K = 2$ для определения выборочного значения \mathcal{P}'_j ($j = 1, \dots, \tau$) можно использовать следующее выражение:

$$\mathcal{P}'_j = q_1 \int_{|d_1|}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy + q_2 \int_{|d_2|}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy,$$

где $|d_1|$ и $|d_2|$ - расстояния от центров многомерных нормальных распределений $\mu^{(1)}$ и $\mu^{(2)}$ до гиперплоскости:

$$(\bar{\mu}^{(1)}, \bar{\mu}^{(2)})' x - \frac{1}{2} (\bar{\mu}^{(1)} \bar{\mu}^{(2)})' (\bar{\mu}^{(1)} + \bar{\mu}^{(2)}) = 0,$$

где $\bar{\mu}^{(1)}$ и $\bar{\mu}^{(2)}$ оценки величин $\mu^{(1)}$ и $\mu^{(2)}$, определяемые на основе обучающей выборки объёма N_j -ой группы. Такой подход к вычислению P' не связан с получением контрольной выборки, что сокращает машинное время проведения эксперимента.

Автор благодарен доктору Ф.М.Н. С.В. Нагаеву и к.т.н. Н.Г. Загоруйко за ряд критических замечаний.

Л И Т Е Р А Т У Р А

I. A.E. Заславский, Н.М. Сычева. Об одной задаче оптимального распознавания образов.—Вычислительные системы. Новосибирск, 1965, вып. 19., стр. 35

2. Г.С. Лбов. Выбор эффективной системы зависимых признаков.—
Вычислительные системы. Новосибирск, 1965. вып. I9,
стр. 2I.
3. Г.С. Лбов. Алгоритмы и программы для распознавания образов
(диагностики заболеваний сердца по баллистокардио-
граммам). Отчет ИМ СО АН СССР, 1966.
4. T.Marill and D.M.Green. On the Effectiveness of Receptors
in Recognition Systems. IREE Trans. on information
theory, v. IF-9, January, 1963, pp II-I7.
5. Г. Фридман. О математическом ожидании ошибки определения
вероятности ошибочной классификации. - Труды Ин-
ститута инженеров по электротехнике и радиотехнике
(русский перевод) . М. , 1965., том 53, № 6, стр.
760.
6. Т. Андерсон. Введение в многомерный статистический анализ.
М, ФМ. 1963.
7. Н.В. Дунин-Барковский, Н.В. Смирнов. Теория вероятностей и
математическая статистика в технике (общая часть).
М. 1955.

ПРИЛОЖЕНИЕ

Поступило в редакцию
24. у. 1966 г.

Таблица I

$\bar{\rho}_x = 0,0099$												
$\rho_x = 4,66$												
$n = 2$				$n = 6$				$n = 10$				
$N = 20$	$N = 50$	$N = 100$	$N = 150$	$N = 20$	$N = 50$	$N = 100$	$N = 150$	$N = 20$	$N = 50$	$N = 100$	$N = 150$	
$\bar{\rho}'_1$	0,0050	0,0060	0,0045	0,0045	0,0090	0,0080	0,0060	0,0065	0,0120	0,0075	0,0100	0,0100
$\bar{\rho}'_2$	0,0085	0,0050	0,0045	0,0045	0,0090	0,0085	0,0055	0,0045	0,0125	0,0120	0,0085	0,0075
$\bar{\rho}'_3$	0,0065	0,0045	0,0045	0,0045	0,0110	0,0065	0,0055	0,0075	0,0100	0,0075	0,0075	0,0100
$\bar{\rho}'_4$	0,0050	0,0055	0,0045	0,0045	0,0065	0,0060	0,0075	0,0060	0,0165	0,0085	0,0085	0,0075
$\bar{\rho}'_5$	0,0050	0,0055	0,0050	0,0045	0,0120	0,0060	0,0065	0,0045	0,0125	0,0110	0,0095	0,0090
$\bar{\rho}'_6$	0,0050	0,0045	0,0045	0,0045	0,0075	0,0075	0,0060	0,0055	0,0115	0,0095	0,0095	0,0085
$\bar{\rho}'_7$	0,0045	0,0045	0,0045	0,0050	0,0080	0,0090	0,0065	0,0070	0,0095	0,0110	0,0075	0,0080
$\bar{\rho}'_8$	0,0055	0,0060	0,0050	0,0045	0,0065	0,0080	0,0065	0,0070	0,0085	0,0095	0,0080	0,0085
$\bar{\rho}'_9$	0,0055	0,0050	0,0045	0,0050	0,0075	0,0065	0,0055	0,0065	0,0105	0,0115	0,0090	0,0090
$\bar{\rho}'_{10}$	0,0055	0,0045	0,0050	0,0045	0,0080	0,0055	0,0065	0,0055	0,0085	0,0085	0,0085	0,0090
$\bar{\rho}'_{11}$	0,0060	0,0050	0,0075	0,0045	0,0080	0,0085	0,0060	0,0070	0,0145	0,0085	0,0095	0,0090
$\bar{\rho}'_{12}$	0,0080	0,0045	0,0045	0,0045	0,0105	0,0055	0,0065	0,0060	0,0100	0,0100	0,0080	0,0100
$\bar{\rho}'_{13}$	0,0050	0,0050	0,0050	0,0045	0,0095	0,0060	0,0065	0,0060	0,0085	0,0075	0,0085	0,0080
$\bar{\rho}'_{14}$	0,0045	0,0045	0,0050	0,0045	0,0065	0,0065	0,0060	0,0060	0,0115	0,0095	0,0095	0,0070
$\bar{\rho}'_{15}$	0,0050	0,0045	0,0045	0,0045	0,0065	0,0060	0,0070	0,0060	0,0165	0,0090	0,0100	0,0085
$\bar{\rho}'_{max}$	0,0085	0,0060	0,0075	0,0050	0,0120	0,0090	0,0075	0,0075	0,0165	0,0120	0,0100	0,0100
$\bar{\rho}'_{min}$	0,0045	0,0045	0,0045	0,0045	0,0065	0,0055	0,0055	0,0045	0,0085	0,0075	0,0075	0,0070
$\Delta \bar{\rho}'$	0,0040	0,0005	0,0030	0,0005	0,0055	0,0085	0,0020	0,0030	0,0080	0,0045	0,0025	0,0030

Таблица 2

$$\bar{\rho}_2 = 0,1003$$

$$\rho_2 = 2,56$$

	$n = 2$				$n = 6$				$n = 10$			
	$N=20$	$N=50$	$N=100$	$N=150$	$N=20$	$N=50$	$N=100$	$N=150$	$N=20$	$N=50$	$N=100$	$N=150$
ρ'_1	0,0890	0,0980	0,0890	0,0800	0,1080	0,1010	0,1005	0,1030	0,1225	0,1185	0,1160	0,1060
ρ'_2	0,0970	0,0910	0,0930	0,0920	0,1055	0,1060	0,0990	0,1020	0,1225	0,1130	0,1075	0,1075
ρ'_3	0,0920	0,0880	0,0905	0,0870	0,1205	0,1023	0,1000	0,0950	0,1215	0,1155	0,1100	0,1095
ρ'_4	0,0885	0,0925	0,0895	0,0890	0,1040	0,1035	0,0975	0,0990	0,1218	0,1110	0,1065	0,1030
ρ'_5	0,0895	0,0930	0,0915	0,0895	0,1110	0,1090	0,1055	0,1065	0,1115	0,1195	0,1120	0,1120
ρ'_6	0,0935	0,0920	0,0900	0,0905	0,0980	0,1000	0,0985	0,0990	0,1315	0,1050	0,1090	0,1120
ρ'_7	0,0920	0,0905	0,0895	0,0940	0,1010	0,1070	0,0985	0,1015	0,1190	0,1130	0,1050	0,1080
ρ'_8	0,0955	0,0905	0,0875	0,0915	0,1065	0,1000	0,0980	0,0990	0,1200	0,1130	0,1175	0,1075
ρ'_9	0,0925	0,0935	0,0900	0,0980	0,1295	0,1010	0,1045	0,1025	0,1260	0,1115	0,1105	0,1095
ρ'_{10}	0,0930	0,0920	0,0925	0,0895	0,1060	0,1015	0,0985	0,1025	0,1325	0,1110	0,1000	0,1070
ρ'_{11}	0,0925	0,0915	0,0925	0,0915	0,1165	0,1000	0,1020	0,0990	0,1230	0,1070	0,1065	0,1090
ρ'_{12}	0,0950	0,0890	0,0980	0,0870	0,1195	0,1040	0,1005	0,1005	0,1160	0,1035	0,1070	0,1120
ρ'_{13}	0,0920	0,0905	0,0925	0,0890	0,0995	0,1030	0,0955	0,0980	0,1330	0,1160	0,1060	0,1120
ρ'_{14}	0,0900	0,0890	0,0875	0,0915	0,1025	0,1065	0,0995	0,0995	0,1160	0,1195	0,1055	0,1095
ρ'_{15}	0,0985	0,0915	0,0925	0,0915	0,1100	0,1010	0,1065	0,1010	0,1385	0,1090	0,1065	0,1160
ρ'^{max}	0,0985	0,0935	0,0930	0,0940	0,1295	0,1090	0,1065	0,1080	0,1385	0,1195	0,1175	0,1160
ρ'^{min}	0,0885	0,0880	0,0875	0,0800	0,0980	0,1000	0,0955	0,0950	0,1115	0,1035	0,1050	0,1080
$\Delta \rho'$	0,0100	0,0055	0,0055	0,0140	0,0315	0,0090	0,0110	0,0080	0,0220	0,0160	0,0125	0,0130

Таблица 3

$$\bar{\rho}_3 = 0,2981$$

$$\rho_3 = 1,06$$

	$n = 2$				$n = 6$				$n = 10$			
	$N=20$	$N=50$	$N=100$	$N=150$	$N=20$	$N=50$	$N=100$	$N=150$	$N=20$	$N=50$	$N=100$	$N=150$
ρ'_1	0,3055	0,3060	0,3075	0,3030	0,3410	0,3160	0,3210	0,3285	0,3820	0,3380	0,3145	0,3105
ρ'_2	0,3055	0,3130	0,3065	0,3105	0,3300	0,3290	0,3065	0,3095	0,3485	0,3480	0,3050	0,3155
ρ'_3	0,3045	0,3075	0,3034	0,3050	0,3385	0,3135	0,3135	0,3095	0,3225	0,3205	0,3110	0,3150
ρ'_4	0,3030	0,3115	0,3085	0,3080	0,3435	0,3260	0,3115	0,3190	0,3530	0,3285	0,3050	0,3240
ρ'_5	0,3060	0,3075	0,3060	0,3065	0,3185	0,3195	0,3175	0,3080	0,3295	0,3170	0,3335	0,3180
ρ'_6	0,3090	0,3075	0,3060	0,3065	0,3105	0,3145	0,3170	0,3145	0,3660	0,3185	0,3155	0,3105
ρ'_7	0,3090	0,3070	0,3030	0,3065	0,3195	0,3245	0,3130	0,3170	0,3090	0,3155	0,3165	0,3034
ρ'_8	0,3075	0,3034	0,3080	0,3055	0,3225	0,3135	0,3190	0,3130	0,3515	0,3265	0,3115	0,2995
ρ'_9	0,3155	0,3145	0,3085	0,3080	0,3360	0,3205	0,3200	0,3165	0,3485	0,3850	0,3125	0,3050
ρ'_{10}	0,3120	0,3090	0,3070	0,3020	0,3160	0,3200	0,3245	0,3225	0,3645	0,3265	0,3070	0,3135
ρ'_{11}	0,3055	0,3150	0,3070	0,3050	0,3420	0,3165	0,3300	0,3140	0,3325	0,3345	0,3145	0,3055
ρ'_{12}	0,3045	0,3070	0,3075	0,3060	0,3235	0,3130	0,3110	0,3090	0,3120	0,3195	0,3230	0,3115
ρ'_{13}	0,3080	0,3055	0,3055	0,3065	0,3175	0,3185	0,3060	0,3155	0,3510	0,3255	0,3220	0,3135
ρ'_{14}	0,3155	0,3034	0,3060	0,3060	0,3135	0,3210	0,3085	0,3125	0,3350	0,3515	0,3110	0,3065
ρ'_{15}	0,3185	0,3070	0,3100	0,3065	0,3405	0,3240	0,3125	0,3205	0,3670	0,3210	0,3140	0,3265
ρ'^{max}	0,3185	0,3150	0,3100	0,3105	0,3435	0,3290	0,3300	0,3285	0,3820	0,3515	0,3335	0,3265
ρ'^{min}	0,3050	0,3034	0,3020	0,3020	0,3105	0,3130	0,3060	0,3080	0,3090	0,3135	0,3050	0,3034
$\Delta \rho'$	0,0155	0,0126	0,0070	0,0085	0,0330	0,0160	0,0240	0,0205	0,0730	0,0380	0,0285	0,0231

Т а б л и ц а 4

$$\mathcal{P}_4 = 0,4483$$

$$\rho_4 = 0,13$$

	$n = 2$				$n = 6$				$n = 10$			
	$N=20$	$N=50$	$N=100$	$N=150$	$N=20$	$N=50$	$N=100$	$N=150$	$N=20$	$N=50$	$N=100$	$N=150$
\mathcal{P}'_1	0,4820	0,4915	0,4745	0,4760	0,4960	0,4885	0,4805	0,4925	0,5070	0,5000	0,4865	0,4975
\mathcal{P}'_2	0,4745	0,4855	0,4750	0,4940	0,4900	0,4915	0,4840	0,4785	0,4935	0,4855	0,4905	0,4985
\mathcal{P}'_3	0,5020	0,4740	0,4735	0,4765	0,4895	0,4855	0,4850	0,4845	0,4895	0,4945	0,4840	0,4865
\mathcal{P}'_4	0,4770	0,4800	0,4755	0,5045	0,5095	0,5095	0,4940	0,4855	0,4995	0,4935	0,4935	0,4885
\mathcal{P}'_5	0,4715	0,4775	0,4785	0,4715	0,4790	0,4885	0,4940	0,4750	0,4915	0,4905	0,5015	0,4925
\mathcal{P}'_6	0,4745	0,4745	0,4705	0,4715	0,4945	0,4870	0,4980	0,4895	0,4945	0,4795	0,4955	0,4820
\mathcal{P}'_7	0,4775	0,4790	0,4785	0,5045	0,4810	0,4990	0,4820	0,4940	0,4775	0,4755	0,4770	0,4880
\mathcal{P}'_8	0,4775	0,4865	0,4840	0,4690	0,5060	0,4805	0,5000	0,4935	0,4965	0,4870	0,5005	0,5050
\mathcal{P}'_9	0,5005	0,4845	0,4765	0,4800	0,5100	0,4810	0,4810	0,4955	0,4890	0,4805	0,4845	0,5015
\mathcal{P}'_{10}	0,4990	0,4755	0,5075	0,4725	0,4855	0,4995	0,5120	0,4890	0,4990	0,4975	0,4810	0,5000
\mathcal{P}'_{11}	0,4855	0,5055	0,4755	0,4925	0,4980	0,5025	0,4810	0,4855	0,4925	0,4850	0,4985	0,5025
\mathcal{P}'_{12}	0,4770	0,4810	0,4760	0,4825	0,5120	0,4905	0,4890	0,4895	0,4925	0,4905	0,4935	0,4810
\mathcal{P}'_{13}	0,5010	0,5010	0,4765	0,4840	0,5060	0,5095	0,4875	0,4795	0,4900	0,5050	0,4825	0,4960
\mathcal{P}'_{14}	0,4860	0,4750	0,4795	0,4790	0,4880	0,4910	0,4805	0,4780	0,4850	0,4945	0,4975	0,5065
\mathcal{P}'_{15}	0,5075	0,4800	0,4750	0,4500	0,4995	0,5010	0,4850	0,4980	0,4950	0,4975	0,4825	0,4900
\mathcal{P}'_{max}	0,5020	0,5010	0,5075	0,5045	0,5120	0,5095	0,5120	0,4955	0,5070	0,5050	0,5015	0,5065
\mathcal{P}'_{min}	0,4715	0,4740	0,4785	0,4690	0,4790	0,4805	0,4810	0,4750	0,4775	0,4755	0,4770	0,4810
$\Delta \mathcal{P}'$	0,0305	0,0270	0,0340	0,0355	0,0380	0,0290	0,0310	0,0205	0,0295	0,0295	0,0245	0,0255