

М.И. Шлезингер, А.В. Бондаренко

## Как формулировать задачи обучения в распознавании образов

Исследованы задачи распознавания образов в ситуации, когда статистическая модель распознаваемого объекта известна лишь частично. Выполнен критический анализ минимаксного подхода к решению таких задач и подхода, основанного на максимально правдоподобном оценивании модели по обучающей выборке. Сформулирована постановка задачи, покрывающая весь спектр ситуаций для обучающих выборок любого объема, от нулевого до бесконечного. Выполнен формальный анализ задач обучения в этой новой постановке и показано ее решение в некоторых простейших случаях.

Pattern recognition problems are considered for a case when a statistical model of an object is not completely known. A minimax approach to solution of such problems is critically analyzed as well as an approach based on the maximal likelihood model estimation with respect to given training multiset. The suggested formulation of the recognition learning problem embraces a whole spectrum of situations for training sets of an arbitrary size: from zero to infinite ones. Main formal properties of the suggested problem formulation are analyzed and its solution in several simplest cases is shown.

Досліджено задачі розпізнавання образів у ситуації, коли статистична модель розпізнаваного об'єкта відома лише частково. Виконано критичний аналіз мінімаксного підходу до розв'язання таких задач та підходу, заснованого на максимально правдоподібному оцінюванні моделі за навчальною вибіркою. Сформульовано постановку задачі, яка покриває весь спектр ситуацій для навчальних вибірок будь-якого об'єму, від нульового до безкінечного. Виконано формальний аналіз задач навчання у цій новій постановці та показано її розв'язання у деяких найпростіших випадках.

**Введение.** В работе исследуется задача распознавания образов в обычной для практики ситуации, когда статистическая модель распознаваемого объекта известна не полностью, а известен лишь класс моделей, которому она принадлежит. Известно, например [1], что задача распознавания допускает разумную, так называемую минимаксную формулировку и в такой, не полностью определенной статистической ситуации. Задача в этой формулировке состоит в построении стратегии, обеспечивающей удовлетворительное качество распознавания при любой модели из заданного класса.

В случае, когда недостающие знания о модели объекта частично восполняются так называемой обучающей выборкой, возникает вопрос, как эти дополнительные сведения использовать для улучшения качества распознавания. Совокупность подходов к решению этого вопроса составляет современную теорию обучения в распознавании образов. Большая часть этих подходов состоит в разделении процесса распознавания на два этапа. Сначала находят модель из заданного класса, которая в определенном смысле наилучшим образом согласо-

ется с обучающей выборкой. Затем с этой наилучшей моделью поступают так, будто она совпадает с действительной. Это отождествление правомерно лишь для обучающих выборок достаточно большой длины. В противном случае возникает масса вопросов, известных как проблема коротких выборок.

Указанные подходы охватывают далеко не всю проблему, а лишь два ее крайних частных случая: в отсутствие обучающей выборки (минимаксный подход) и при неограниченно большой ее длине (существующие методы обучения). Дадим такую формулировку задачи обучения в распознавании, которая охватывает весь спектр ситуаций для обучающих выборок любой длины: от нулевой до бесконечной.

### Формулировка основных понятий

Будем говорить об объекте, который характеризуется двумя случайными параметрами  $x$  и  $k$ . Эти параметры принимают значения из двух конечных множеств,  $X$  и  $K$  соответственно. Параметр  $x$  будем называть *наблюдаемым признаком* объекта, или *наблюдением*, а параметр  $k$  – его *скрытым*, непосредственно *не наблюдаемым состоянием*. Также будем говорить о

конечном множестве  $M$  моделей и об  $t \in M$ , как о некоторой модели из этого множества. Пусть функция  $p: X \times K \times M \rightarrow R$  такова, что ее значение  $p(x, k; m)$ ,  $x \in X$ ,  $k \in K$ ,  $m \in M$ , есть совместная вероятность наблюдения  $x$  и состояния  $k$  в модели  $m$ .

Множество  $\{(x, k) | x \in X, k \in K\}$  обозначим  $\Lambda$ , а через  $\Lambda^n$  обозначим множество обучающих выборок длины  $n$ ,  $n = 1, \dots$ . Одна отдельная выборка  $L^n \in \Lambda^n$  имеет вид

$$L^n = (x_1, k_1; x_2, k_2; \dots, x_n, k_n).$$

Введем также обозначение  $L^0$  для обучающей выборки нулевой длины и  $\Lambda^0$  для непустого множества, содержащего единственную выборку  $L^0$ .

Примем, что на множестве  $\Lambda^n$  задано распределение вероятностей, зависящее от модели  $m$ , так что

$$p(L^n; m) = \prod_{i=1}^n p(x_i, k_i; m), \quad L^n \in \Lambda^n, \quad m \in M.$$

При этом считается, что  $p(L^0; m) = 1$ .

Функцию  $q: X \times K \rightarrow R$  будем называть *стратегией распознавания*. Ее значение  $q(k/x)$  понимается как условная вероятность принятия решения, что объект находится в состоянии  $k$ , при условии, что на вход стратегии подано наблюдение  $x$ . На данной стадии изложения введение такой рандомизированной стратегии распознавания может показаться ненужным усложнением в сравнении с детерминированной стратегией вида  $X \rightarrow K$ . Однако в дальнейшем определение стратегии именно в рандомизированном виде существенно упростит доказательство некоторых утверждений.

Множество  $R^{X \times K}$  всех возможных стратегий вида  $q: X \times K \rightarrow R$  обозначим  $Q$ . При этом, естественно, предполагается, что  $Q$  содержит только те функции  $q$ , для которых

$$\sum_{k \in K} q(k/x) = 1, \quad x \in X, \quad q(k/x) \geq 0, \quad x \in X, \quad k \in K.$$

Под стратегией обучения, или алгоритмом обучения, будем понимать функцию, которая

для каждой обучающей выборки  $L^n$  указывает стратегию распознавания из множества  $Q$ . Обозначим  $q^n: \Lambda^n \rightarrow Q$  стратегию обучения по обучающей выборке длины  $n$ . Для каждой обучающей выборки  $L^n \in \Lambda^n$  она определит стратегию  $q^n(L^n): X \times K \rightarrow R$  распознавания. Эта стратегия распознавания, примененная к некоторому конкретному значению  $x$  признака, выдает случайное решение  $k$  о состоянии объекта с вероятностью  $q^n(L^n)(k/x)$ . Для каждой стратегии  $q^n: \Lambda^n \rightarrow Q$  определим рандомизированную стратегию  $\tilde{q}^n: L^n \times X \times K \rightarrow R$ , так что

$$\tilde{q}^n(k/x, L^n) = q^n(L^n)(k/x),$$

и назовем ее *стратегией обучения-распознавания*. Стратегии  $\tilde{q}^n$  и  $q^n$  – формально различные функции, так как первая имеет формат  $\Lambda^n \times X \times K \rightarrow R$ , а вторая –  $\Lambda^n \rightarrow Q$ . Однако их содержательный смысл один и тот же. Стратегия обучения  $q^n: \Lambda^n \rightarrow Q$  указывает, как стратегия распознавания зависит от обучающей выборки, а стратегия обучения-распознавания  $\tilde{q}^n: L^n \times X \times K \rightarrow R$  указывает, как решение о состоянии  $k$  объекта по признаку  $x$  должно зависеть от обучающей выборки  $L^n$ . Поэтому в дальнейшем будем использовать как понятие стратегии обучения  $q^n$ , так и равноценное ему понятие стратегии  $\tilde{q}^n$  обучения-распознавания.

Пусть  $w: K \times K \rightarrow R$  – функция потерь со значениями  $w(k, k')$ , обозначающими потери в случае, когда об объекте, находящемся в состоянии  $k$ , принимается решение, что он находится в состоянии  $k'$ . Риском стратегии  $q: X \times K \rightarrow R$  на модели  $m$  называется математическое ожидание штрафа, т.е. число

$$R(q, m) = \sum_{\substack{k \in K \\ x \in X}} \sum_{k' \in K} q(k'/x) \cdot p(x, k; m) \cdot w(k, k').$$

Для каждой модели  $m \in M$  определим риск  $R^B(m) = \min_{q \in Q} R(q, m)$ , который назовем байесовским риском для модели  $m$ , и стратегию

$q^B(m) = \arg \min_{q \in Q} R(q, m)$ , которую назовем байесовской стратегией для модели  $m$ . Величина  $R(q^n(L^n), m)$  – это риск, который на модели  $m$  достигается распознающей системой после ее обучения на обучающей выборке  $L^n$ . В силу случайного характера выборки  $L^n$  риск  $R(q^n(L^n), m)$  также будет случайным. Его математическое ожидание по всем выборкам  $L^n \in \Lambda^n$  обозначим

$$R^n(q^n, m) = \sum_{L^n \in \Lambda^n} p(L^n; m) \cdot R(q^n(L^n), m) = \sum_{L^n \in \Lambda^n} \prod_{i=1}^n p(x_i, k_i; m) \cdot R(q^n(L^n), m). \quad (1)$$

Задача обучения в распознавании образов в ее неформальном понимании состоит в том, чтобы для заданного числа  $n$  и заданных функций  $p: X \times K \times M \rightarrow R$  и  $w: K \times K \rightarrow R$  построить удовлетворительную в определенном смысле стратегию  $q^n$  обучения или, что то же самое, стратегию  $\tilde{q}^n$  обучения-распознавания. Однако качество искомой стратегии определяется не одним числом, а совокупностью чисел  $R^n(q^n, m)$ ,  $m \in M$ .

Современные исследования проблемы обучения состоят совсем не в конкретизации приведенной формулировки и получении процедур обучения-распознавания как решения конкретно сформулированной задачи, а в фиксировании некоторых общепризнанных процедур обучения и последующем анализе ситуаций, когда они оказываются особенно плохими. Укажем наиболее известные такие процедуры.

**Известные подходы к распознаванию при не полностью определенной модели распознаваемого объекта**

#### Минимаксная стратегия

Потребуем, чтобы стратегия распознавания  $q$  обеспечивала удовлетворительное качество распознавания для любой модели  $m$  из заданного класса  $M$ , т.е. чтобы

$$\{R(q, m) \leq c, m \in M\}. \quad (2)$$

При малых значениях величины  $c$  требования (2) могут оказаться противоречивыми. Поэтому сформулируем задачу отыскания минимального значения  $c$ , при котором требования (2) еще выполнимы. Пусть это минимальное значение равно  $c^*$ . Тогда стратегию  $q^*$ , удовлетворяющую требованиям  $\{R(q^*, m) \leq c^*, m \in M\}$

назовем минимаксной стратегией распознавания на множестве  $M$  моделей. Поскольку задача построения минимаксной стратегии сформулирована для случая отсутствия обучающей выборки, для некоторых классов моделей гарантированный уровень  $c^*$  может оказаться слишком плохим. Само по себе это не является недостатком минимаксных стратегий, потому что для некоторых моделей даже байесовский риск может оказаться неприемлемо большим. Существенный изъян минимаксного подхода состоит в том, что минимаксные стратегии распознавания без обучения не допускают разумного обобщения на минимаксные стратегии обучения. А именно, наилучший в минимаксном смысле алгоритм обучения – это алгоритм, который обучающую выборку игнорирует. Далее подробно рассмотрим этот существенный недостаток минимаксного подхода.

#### Наиболее правдоподобное оценивание модели

Пусть для каждой обучающей выборки  $L^n = (x_1, k_1; x_2, k_2; \dots, x_n, k_n)$  определена наиболее правдоподобная модель

$$m^*(L^n) = \arg \max_{m \in M} \prod_{i=1}^n p(x_i, k_i; m),$$

а для каждой модели  $m \in M$  определена байесовская стратегия  $q^B(m) = \arg \min_{q \in Q} R(q, m)$ .

Стратегия обучения  $q^n: \Lambda^n \rightarrow Q$  определена так, что

$$q^n(L^n) = q^B(m^*(L^n)), L^n \in \Lambda^n. \quad (3)$$

Стратегии такого вида глубоко укоренились в теории и практике распознавания. Именно в силу того, что они хорошо исследованы, об-

щеизвестны их недостатки, обозначаемые кратко, как проблема коротких выборок. Покажем их на простом примере в следующем разделе. Несмотря на эти известные недостатки, стратегия обучения (3) показала свою плодотворность в таких практических задачах, в которых возможны достаточно длинные обучающие выборки. Следует отметить также, что для многих классов моделей как отыскание наиболее правдоподобной модели  $m^*(L^n)$ , так и построение байесовской стратегии  $q^B(m)$  оказывается далеко не тривиальной задачей, представляющей самостоятельный интерес.

### Минимизация эмпирического риска

Пусть  $Q(M) = \{q^B(m) \mid m \in M\} \subset Q$  – множество байесовских стратегий для моделей из класса  $M$ , а  $R^3(q, L^n) = \sum_{i=1}^n \sum_{k \in K} q(k/x_i) \cdot W(k_i, k)$  – так называемый эмпирический риск стратегии  $q: X \times K \rightarrow R$  на обучающей выборке  $L^n = (x_1, k_1; x_2, k_2; \dots, x_n, k_n)$ . Стратегия  $q^n: L^n \rightarrow Q$  определена как

$$q^n(L^n) = \arg \min_{q \in Q(M)} R^3(q, L^n), L^n \in \Lambda^n.$$

Наряду с методом наибольшего правдоподобия минимизация эмпирического риска является одним из наиболее популярных методов обучения. Критический анализ этого подхода к обучению в распознавании выполнен в [1].

### Иллюстрация введенных понятий и известных подходов

Пусть  $X$  – множество вещественных чисел,  $K = \{1, 2\}$  а  $p(x/k)$ ,  $k = 1, 2$ , – распределение плотности условной вероятности, определенное как

$$p(x/k) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_k)^2}, k = 1, 2, \mu_1 = 1, \mu_2 = -1.$$

Пусть  $p_k$ ,  $k = 1, 2$ , – априорная вероятность состояния  $k$ . Вероятности  $p_1$  и  $p_2$  неизвестны, и пара  $(p_1, p_2)$  – неизвестная модель  $m$  объекта, так что  $p(x, k; m) = p_k \cdot p(x/k)$ .

Таким образом, в данном примере статистическая модель распознаваемого объекта почти полностью известна. Неизвестны только апри-

орные вероятности пребывания объекта в первом или втором состояниях. Этот специально упрощенный пример предназначен для предельно ясной иллюстрации введенных формальных понятий, недостатков метода наибольшего правдоподобия и сущности предлагаемого далее подхода к обучению в распознавании.

Пусть потери равны

$$w(k, k') = 1, \text{ если } k \neq k', w(k, k') = 0, \text{ если } k = k'.$$

Байесовская стратегия  $q: X \times K \rightarrow R$  для модели  $m = (p_1, p_2)$  имеет вид

$$q(k = 1/x) = 1, \text{ если } p_1 \cdot p(x/1) > p_2 \cdot p(x/2),$$

$$q(k = 2/x) = 1, \text{ если } p_1 \cdot p(x/1) < p_2 \cdot p(x/2).$$

На рис. 1 показана зависимость риска различных стратегий от модели  $m$ , в данном случае, от вероятности  $p_1$  первого состояния. Риск  $R^B(m)$  байесовской стратегии представлен сплошной линией.

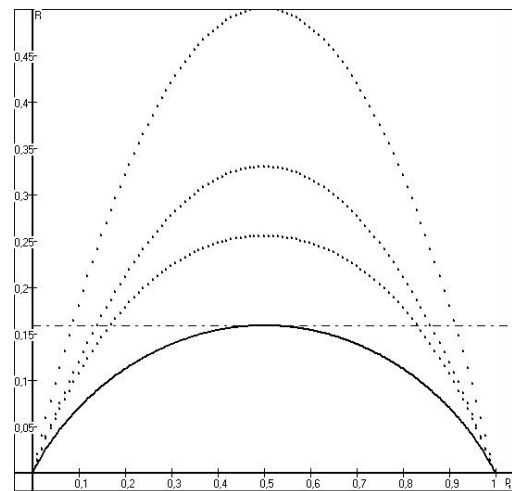


Рис. 1. Зависимость вероятности ошибочного распознавания от априорных вероятностей состояний для различных стратегий

Минимаксная стратегия  $q'(x): X \times K \rightarrow R$  имеет вид

$$q'(k = 1/x) = 1, \text{ если } x > 0,$$

$$q'(k = 2/x) = 1, \text{ если } x < 0.$$

Риск  $R(q', m)$  для этой стратегии показан в виде штрихпунктирной линии.

Пусть  $L^n$  – обучающая выборка длины  $n > 0$ . Наиболее правдоподобная модель  $m^*(L^n)$  или,

что то же самое, наиболее правдоподобные значения  $p_1^*, p_2^*$  априорных вероятностей равны  $m^*(L^n) = (p_1^*, p_2^*) = \left(\frac{n_1}{n}, \frac{n_2}{n}\right)$ , где  $n_k$ ,  $k=1,2$ , – количество элементов  $i$  в выборке  $L^n$ , для которых  $k_i = k$ . Стратегия  $\tilde{q}^n : \Lambda^n \times X \times K \rightarrow R$  в этом случае имеет вид

$$\tilde{q}^n(k=1/L^n, x) = 1, \text{ если } n_1 \cdot p(x/1) > n_2 \cdot p(x/2),$$

$$\tilde{q}^n(k=2/L^n, x) = 1, \text{ если } n_1 \cdot p(x/1) < n_2 \cdot p(x/2).$$

Риск  $R^n(\tilde{q}^n, m)$  этой стратегии зависит от длины  $n$  выборки и для значений  $n$ , равных 1, 2 или 3, приведен на рис. 1 в виде пунктирной линии. Очевидно, что при больших значениях длины выборки риск  $R^n(\tilde{q}^n, m)$  становится неотличимым от байесовского риска  $R^B(m)$ . Однако, как видно из рис. 1, при значениях длины  $n$ , равных 1, 2 или 3, риск  $R^n(\tilde{q}^n, m)$  заметно отличается от байесовского. Существенно, что риск  $R^n(\tilde{q}^n, m)$  отличается в худшую сторону и от риска минимаксной стратегии, которая обучающую выборку просто игнорирует. Следовательно, метод наибольшего правдоподобия неэффективно использует информацию, содержащуюся в обучающейся выборке, потому что ее использование приводит к результатам, худшим, чем ее игнорирование.

В следующем разделе исследуется вопрос о том, можно ли построить минимаксные стратегии обучения-распознавания так, чтобы использование коротких выборок пусть незначительно, но улучшало бы результаты в сравнении с теми, которые получаются в отсутствие обучающей выборки. На этот вопрос будет получен отрицательный ответ.

### Минимаксные стратегии обучения-распознавания

Сформулируем задачу отыскания такого минимального значения  $c$ , при котором еще не противоречивы требования  $\{R^n(\tilde{q}^n, m) \leq c, m \in M\}$ , к стратегии  $\tilde{q}^n$  обучения-распознавания.

Учитывая, что для  $n=0$  стратегия  $q^0 : L^0 \times X \times K \rightarrow R$  есть стратегия вида  $X \times K \rightarrow R$  и

$$R^0(q^0, m) = R(q, m) = \sum_{\substack{k' \in K \\ x \in X}} q(k'/x) \sum_{k \in K} p(x, k; m) \cdot w(k, k'),$$

запишем эту задачу для  $n=0$  в виде следующей пары двойственных задач линейного программирования.

Прямая задача:

$$\min c \quad (4)$$

$$\tau(m) \left| \sum_{\substack{k' \in K \\ x \in X}} q(k'/x) \sum_{k \in K} p(x, k; m) \cdot w(k, k') \leq c, \quad (5)$$

$$t(x) \left| \begin{array}{l} m \in M; \\ \sum_{k' \in K} q(k'/x) = 1, \quad x \in X; \\ q(k'/x) \geq 0, \quad x \in X; \quad k' \in K. \end{array} \quad (6)$$

Двойственная задача

$$\max \sum_{x \in X} t(x) \quad (7)$$

$$q(k'/x) \left| t(x) \leq \sum_{k \in K} \left( \sum_{m \in M} \tau(m) \cdot p(x, k; m) \right) \times \quad (8)$$

$$c \left| \begin{array}{l} \times w(k, k'), \quad k' \in K, \quad x \in X; \\ \sum_{m \in M} \tau(m) = 1; \quad \tau(m) \geq 0, \quad x \in X. \end{array} \quad (9)$$

В приведенных формулах слева от каждого ограничения прямой задачи записана соответствующая ей двойственная переменная; слева от каждого ограничения двойственной задачи записана соответствующая ей переменная в прямой задаче.

Функцию вида  $\tau : M \rightarrow R$  назовем весовой функцией со значениями  $\tau(m)$ , которые назовем весами. Множество весовых функций, удовлетворяющих ограничениям  $\sum_{m \in M} \tau(m) = 1, \tau(m) \geq 0, m \in M$ , обозначим  $T$ . Определим модель  $m(\tau)$  так, что совместная вероятность  $x$  и  $k$  в модели  $m(\tau)$  определяется выражением

$$p(x, k; m(\tau)) = \sum_{m \in M} \tau(m) \cdot p(x, k; m). \quad (10)$$

Модели этого вида назовем *смешанными* моделями в отличие от моделей из первоначального класса  $M$ , которые будем называть *чистыми*. Ранее введенные стратегии  $q^B(m)$ ,  $m \in M$ , будем называть чистыми в отличие от стратегий  $\arg \min_{q \in Q} R(q, m(\tau))$ ,  $\tau \in T$ , которые будем называть смешанными.

**Лемма 1.** Пусть  $q^*$  и  $c^*$  – решение прямой задачи (4). В таком случае существует такой набор  $\tau^*(m)$ ,  $m \in M$ , весов моделей, что

$$\max_{m \in M} R(q^*, m) = R^B(m(\tau^*)). \quad (11)$$

**Доказательство.** Пусть  $t^*(x)$  и  $\tau^*(m)$  – решение двойственной задачи (7). В силу известной теоремы двойственности [4]

$$\sum_{x \in X} t^*(x) = c^*. \quad (12)$$

Выражение в левой части (5) есть риск  $R(q^*, m)$ , а  $c^*$  – минимальное число, удовлетворяющее условиям  $R(q^*, m) \leq c^*$ ,  $m \in M$ .

Следовательно,

$$c^* = \max_{m \in M} R(q^*, m). \quad (13)$$

Учитывая, что по определению (10)

$$\sum_{m \in M} \tau(m) \cdot p(x, k; m) = p(x, k; m(\tau)),$$

запишем ограничение (8) в виде

$$t(x) \leq \sum_{k \in K} p(x, k; m(\tau)) \cdot w(k, k'), \quad k' \in K.$$

Число  $t^*(x)$  есть максимальное число, удовлетворяющее условиям

$$t^*(x) \leq \sum_{k \in K} p(x, k; m(\tau^*)) \cdot w(k, k'), \quad k' \in K.$$

Следовательно,

$$t^*(x) = \min_{k' \in K} \sum_{k \in K} p(x, k; m(\tau^*)) \cdot w(k, k'),$$

что в свою очередь обозначает, что

$$\sum_{x \in X} t^*(x) = R^B(m(\tau^*)). \quad (14)$$

Подставляя (13) и (14) в (12), получаем (11). Таким образом, набор весов  $\tau^*(m)$ , существование которого требуется доказать, – это веса

$\tau^*(m)$ , являющиеся решением двойственной задачи (7). **Лемма доказана.**

**Лемма 2.** Для любой стратегии  $q^n$  обучения и любой модели  $m \in M$  справедливо неравенство

$$R^n(q^n, m) \geq R^B(m).$$

**Доказательство.** По определению байесовской стратегии неравенство  $R(q^n(L^n), m) \geq R^B(m)$  справедливо для любой выборки  $L^n \in \Lambda^n$ . Следовательно, справедливо и неравенство

$$\sum_{L^n \in \Lambda^n} p(L^n; m) \cdot R(q^n(L^n), m) \geq R^B(m).$$

Левая часть этого неравенства есть  $R^n(q^n, m)$  по определению (1). **Лемма доказана.**

**Теорема 1.** Если для класса  $M$  моделей решение  $q^0 : X \times K \rightarrow R$  минимаксной задачи распознавания (4) есть чистая стратегия, то для любого  $n = 1, 2, \dots$  и любой стратегии  $q^n : L^n \rightarrow Q$  выполняется неравенство

$$\max_{m \in M} R^n(q^n, m) \geq \max_{m \in M} R(q^0, m).$$

**Доказательство.** Пусть  $m^0 = \arg \max_{m \in M} R(q^0, m)$ .

$$\begin{aligned} \text{В таком случае } \max_{m \in M} R^n(q^n, m) &\geq R^n(q^n, m^0) \geq \\ &\geq R^B(m^0) = \max_{m \in M} R(q^0, m). \end{aligned}$$

Первое неравенство в этой цепочке очевидно. Второе – справедливо в силу леммы 2. В силу предположения о том, что стратегия  $q^0$  – чистая, и в силу леммы 1 справедливо равенство в третьем звене цепочки. **Теорема доказана.**

Поскольку данная теорема представляет достаточно существенный отрицательный результат, рассмотрим несколько примеров классов моделей, где минимаксная стратегия распознавания оказывается чистой.

1. Минимаксная стратегия оказывается чистой, конечно, в том случае, когда множество смешанных моделей совпадает с множеством чистых моделей. Эта ситуация, в частности, имеет место в приведенном примере, когда модель распознаваемого объекта почти полностью известна, а неизвестны только априорные вероятности состояний.

2. Пусть  $K \in \{1, 2\}$  – множество состояний объекта, а  $x$  – многомерная гауссова случайная величина с единичной ковариационной матрицей. Условное математическое ожидание этого вектора при условии, что объект находится в состоянии  $k$ , есть  $\mu_k$ . Математические ожидания  $\mu_1$  и  $\mu_2$  неизвестны. Известны, однако, два выпуклых множества  $M_1$  и  $M_2$ , которым они принадлежат. Множество чистых моделей в этом случае есть  $M_1 \times M_2$  и оно, конечно, не равно множеству смешанных моделей. Тем не менее минимаксная стратегия распознавания – чистая.

3. Для теории и практики распознавания обычна ситуация, когда функции  $p_k : X \rightarrow R$ ,  $k \in K$ , со значениями  $p(x/k)$  неизвестны, но известно лишь, что все они входят в заданный класс  $P$  функций, один и тот же для всех  $k \in K$ . В этом случае минимаксная стратегия распознавания также чистая. Эта наилучшая в минимаксном смысле стратегия оказывается совсем плохой. Она гарантирует лишь, что вероятность ошибочного распознавания не превосходит 50% для любой модели. Однако в силу того, что эта стратегия чистая, вступает в силу вывод доказанной теоремы, что никакая стратегия обучения этот показатель не может улучшить.

Исходя из указанного существенного изъяна минимаксного критерия обучения-распознавания и общеизвестных недостатков максимально правдоподобного оценивания модели, приходим к выводу, что проблема распознавания при не полностью известной статистической модели исследована в весьма незначительной части. Известные и широко применяемые методы охватывают лишь два крайних ее случая: когда имеется обучающая выборка неограниченно большой длины и когда обучающей выборки нет вообще. Желательно сформулировать задачу обучения-распознавания так, чтобы риск распознавания при короткой обучающей выборке оказывался хоть ненамного, но лучше, чем риск распознавания без учета этой выборки. Кроме того, необходимо, чтобы при неограниченном увеличении длины вы-

борки качество распознавания становилось неотличимым от байесовского.

### Общий вид стратегий обучения-распознавания

Задача обучения состоит в выборе стратегии  $\tilde{q}^n$ , оптимальной в определенном смысле, из множества всех возможных стратегий вида  $\tilde{q}^n : L^n \times X \times K \rightarrow R$ . Эта задача допускает множество точных формулировок в зависимости от того, как определено качество стратегии, которое следует оптимизировать. При этом следует учесть, что совокупность  $R^n(\tilde{q}^n, m)$ ,  $m \in M$ , рисков позволяет задать на множестве всех возможных стратегий вида  $\tilde{q}^n : L^n \times X \times K \rightarrow R$  естественное отношение частичной упорядоченности. А именно, стратегия  $\tilde{q}_1^n$  не хуже стратегии  $\tilde{q}_2^n$  на классе моделей  $M$ , если для всех моделей  $m \in M$  выполняется неравенство  $R^n(\tilde{q}_1^n, m) \leq R^n(\tilde{q}_2^n, m)$ , а стратегия  $\tilde{q}_2^n$  во всех отношениях хуже стратегии  $\tilde{q}_1^n$ , если указанные неравенства выполняются строго. Формулировка задачи обучения должна быть разумной в том смысле, что ее решение не должно оказаться стратегией, которая во всех отношениях хуже некоторой другой стратегии. На основании этих соображений можно априори сузить множество стратегий обучения-распознавания еще до формулировки задачи обучения, требуя лишь, чтобы эта формулировка была разумной в указанном выше смысле.

Для заданной весовой функции  $\tau$  определим стратегию обучения-распознавания так, что  $\tilde{q}^n(k' / L^n, x) = 1$ , если для всех  $k'' \neq k'$  выполняется неравенство

$$\sum_{k \in K} \left( \sum_{m \in M} \tau(m) \cdot p(L^n; m) \cdot p(x, k; m) \right) \cdot w(k, k') < \sum_{k \in K} \left( \sum_{m \in M} \tau(m) \cdot p(L^n; m) \cdot p(x, k; m) \right) \cdot w(k, k'').$$

Стратегию такого вида назовем байесовской для весовой функции  $\tau$ . Иными словами, при обучающей выборке  $L^n$  и наблюдении  $x$  байе-

совская стратегия принимает решение в пользу состояния  $k^*$ ,

$$k^* = \arg \min_{k' \in K} \sum_{k \in K} \left( \sum_{m \in M} \tau(m) \cdot p(L^n; m) \cdot p(x, k; m) \right) \times w(k, k').$$

При фиксированном классе  $M$  моделей байесовская стратегия может быть той или иной в зависимости от весовой функции  $\tau$ . Докажем, что решение любой разумно поставленной задачи обучения – это байесовская стратегия для некоторой весовой функции  $\tau: M \rightarrow R$ .

**Теорема 2.** Для любой стратегии  $\tilde{q}^{*n}: L^n \times X \times K \rightarrow R$  обучения-распознавания существует байесовская стратегия  $\tilde{q}^n: L^n \times X \times K \rightarrow R$ , для которой неравенство  $R^n(\tilde{q}^n, m) \leq R^n(\tilde{q}^{*n}, m)$  выполняется для любой модели  $m \in M$ .

**Доказательство.** Представим искомую стратегию  $\tilde{q}^n$  как решение следующей задачи линейного программирования:

$$\begin{aligned} & \min c \\ & \tau(m) \left| \begin{aligned} & c - \sum_{L^n \in \Lambda^n} p(L^n; m) \sum_{x \in X} \sum_{k' \in K} \tilde{q}^n(k' / L^n, x) \times \\ & \times \sum_{k \in K} p(x, k; m) \cdot w(k, k') \geq -R^n(\tilde{q}^{*n}, m), \quad (15) \\ & m \in M; \end{aligned} \right. \\ & t(L^n, x) \left| \begin{aligned} & \sum_{k' \in K} \tilde{q}^n(k' / L^n, x) = 1, x \in X, L^n \in \Lambda^n, \quad (16) \\ & \tilde{q}^n(k' / L^n, x) \geq 0, k' \in K, L^n \in \Lambda^n, x \in X, \end{aligned} \right. \end{aligned}$$

которой соответствует двойственная задача

$$\begin{aligned} & \max \left[ \sum_{x \in X} \sum_{L^n \in \Lambda^n} t(L^n, x) - \sum_{m \in M} \tau(m) \cdot R^n(\tilde{q}^{*n}, m) \right] \\ & \tilde{q}^n(k' / L^n, x) \left| \begin{aligned} & t(L^n, x) \leq \sum_{k \in K} \left( \sum_{m \in M} \tau(m) \times \right. \\ & \times p(L^n; m) \cdot p(x, k; m) \right) \times \\ & \times w(k, k'), k' \in K, L^n \in \Lambda^n, x \in X; \quad (17) \\ & c \left| \begin{aligned} & \sum_{m \in M} \tau(m) = 1; \tau(m) \geq 0, m \in M. \end{aligned} \right. \end{aligned} \right. \end{aligned}$$

Как и ранее, слева от уравнений-ограничений задач записаны соответствующие им двойственные переменные.

Система линейных ограничений прямой задачи не противоречива, так как ей удовлетворяют значение  $c = 0$  и значения  $\tilde{q}^n(k' / L^n, x) = \tilde{q}^{*n}(k' / L^n, x)$ ,  $k' \in K$ ,  $L^n \in \Lambda^n$ ,  $x \in X$ . Нетрудно также увидеть, что величина  $c$  ограничена снизу на множестве решений этой системы. Отсюда, в силу известных теорем двойственности, следует, что ограничения в задаче (17) также непротиворечивы, а целевая функция в этой задаче ограничена сверху.

Пусть  $\tilde{q}_0^n: L^n \times K \times X \rightarrow R$ ,  $c_0$  – решение прямой задачи, а  $t_0(L^n, x)$ ,  $L^n \in \Lambda^n$ ,  $x \in X$ , и  $\tau_0(m)$ ,  $m \in M$  – решение двойственной. Пусть для некоторой тройки  $k_0 \in K$ ,  $L_0^n \in \Lambda^n$ ,  $x_0 \in X$  неравенство

$$\begin{aligned} & \sum_{k \in K} \left( \sum_{m \in M} \tau_0(m) \cdot p(L_0^n; m) \cdot p(x_0, k; m) \right) \cdot w(k, k_0) < \\ & < \sum_{k \in K} \left( \sum_{m \in M} \tau_0(m) \cdot p(L_0^n; m) \cdot p(x_0, k; m) \right) \cdot w(k, k') \end{aligned}$$

выполняется для всех  $k' \neq k_0$ . Это значит, что неравенство (17) выполняется строго для данных  $L_0^n, x_0$  и всех  $k' \neq k_0$ . Отсюда, в силу известной теоремы двойственности, равенство  $q_0^n(k' / L_0^n, x_0) = 0$  выполняется для всех  $k' \neq k_0$ . В свою очередь, в силу (16) отсюда следует, что  $q_0^n(k_0 / L_0^n, x_0) = 1$ .

Таким образом, стратегия  $q_0^n$  является байесовской для весовой функции  $\tau_0$ .

Поскольку значение  $c = 0$  достигается на множестве допустимых решений системы ограничений прямой задачи, то  $c_0 \leq 0$ . Отсюда, в силу неравенств (15), следует, что неравенство

$$\begin{aligned} & \sum_{L^n \in \Lambda^n} p(L^n; m) \sum_{x \in X} \sum_{k' \in K} \tilde{q}_0^n(k' / L^n, x) \sum_{k \in K} p(x, k; m) \times \\ & \times w(k, k') \leq R^n(q^{*n}, m) \end{aligned}$$



выполняется для всех  $m \in M$ . Выражение в левой части этого неравенства есть  $R^n(\tilde{q}_0^n, m)$  и  $R^n(\tilde{q}_0^n, m) \leq R^n(\tilde{q}^{*n}, m)$ ,  $m \in M$ . **Теорема доказана.**

Таким образом, решение любой разумно поставленной задачи обучения существенно отличается от общепринятых стратегий, указанных ранее. В решении разумно поставленной задачи отсутствует такой этап, как выбор одной-единственной модели на основании обучающей выборки, пусть даже наиболее правдоподобной модели или модели с минимальным эмпирическим риском. В распознавании после обучения учитываются все модели из заданного класса, но с различными весами в зависимости от того, насколько та или иная модель согласуется с обучающей выборкой.

**Одна из возможных формулировок задачи обучения-распознавания и ее формальные свойства**

Риск  $R^n(q^n, m)$  зависит не только от стратегии  $q^n$ , но и от модели  $m$ . Для некоторых моделей  $m$  даже байесовский риск  $R^B(m)$  может оказаться неприемлемо большим. Поэтому характеристикой собственно стратегии  $q^n$  естественно считать не столько величину  $R^n(q^n, m)$ , сколько ее отклонение  $R^n(q^n, m) - R^B(m)$  от байесовского риска  $R^B(m)$ , который характеризует собственно модель  $m$ , вне ее связи со стратегией  $q^n$ . Именно исходя из разности  $R^n(q^n, m) - R^B(m)$  исследуются отдельные процедуры обучения в работах [2, 3] и других. Качеством же стратегии  $q^n$  безотносительно модели  $m$  примем величину  $\max_m [R^n(q^n, m) - R^B(m)]$  и сформулируем задачу обучения, как поиск стратегии  $q^{*n}$  с наилучшим качеством, т.е.

$$q^{*n} = \arg \min_{q^n} \max_m [R^n(q^n, m) - R^B(m)]. \quad (18)$$

**Теорема 3.** Искомая стратегия  $\tilde{q}^{*n}$  есть байесовская стратегия для смешанной модели с весовой функцией

$$\tau^* = \arg \max_{\tau \in T} \left[ \min_{q^n} R^n(q^n, m(\tau)) - \sum_{\tau \in T} \tau(m) R^B(m) \right].$$

**Доказательство.** Запишем искомую стратегию как решение следующей задачи линейного программирования.

Прямая задача

$$\begin{aligned} & \min c \\ \tau(m) \quad & \left| \begin{aligned} & \sum_{L^n \in \Lambda^n} \sum_{x \in X} \sum_{k' \in K} \tilde{q}^n(k' / L^n, x) \times \\ & \times p(L^n; m) \sum_{k \in K} p(x, k; m) \times \\ & \times w(k, k') - R^B(m) \leq c, \quad m \in M; \end{aligned} \right. \quad (19) \\ t(L^n, x) \quad & \left| \begin{aligned} & \sum_{k' \in K} \tilde{q}^n(k' / L^n, x) = 1; \quad L^n \in \Lambda^n, x \in X; \\ & \tilde{q}^n(k' / L^n, x) \geq 0, \quad k' \in K, L^n \in \Lambda^n, x \in X. \end{aligned} \right. \end{aligned}$$

Переменными в этой задаче являются числа  $\tilde{q}^n(k' / L^n, x)$ ,  $k' \in K$ ,  $L^n \in \Lambda^n$ ,  $x \in X$ , и число  $c$ . Задаче соответствует двойственная задача с двойственными переменными  $\tau(m)$ ,  $m \in M$ , и  $t(L^n, x)$ ,  $L^n \in \Lambda^n$ ,  $x \in X$ .

Двойственная задача

$$\begin{aligned} & \max \left[ \sum_{x \in X} \sum_{L^n \in \Lambda^n} t(L^n, x) - \sum_{m \in M} \tau(m) \cdot R^B(m) \right] \\ \tilde{q}^n(k' / L^n, x) \quad & \left| \begin{aligned} & t(L^n, x) \leq \sum_{k \in K} \left( \sum_{m \in M} \tau(m) \times \right. \\ & \times p(L^n; m) \times p(x, k; m) \right) \times w(k, k'), \\ & k' \in K, L^n \in \Lambda^n, m \in M; \\ c \quad & \left. \sum_{m \in M} \tau(m) = 1; \tau(m) \geq 0, m \in M. \right. \end{aligned} \right. \end{aligned}$$

Слева от ограничений прямой (или двойственной) задачи записаны соответствующие им переменные двойственной (или прямой) задачи.

Пусть  $c^*$  и  $q^{*n}$  – решение прямой задачи, а  $t^*(L^n, x)$ ,  $L^n \in \Lambda^n$ ,  $x \in X$  и  $\tau^*(m)$ ,  $m \in M$ , – решение двойственной. Пусть  $L_0^n$ ,  $x_0$ ,  $k_0$  – такие значения, что для любого  $k' \neq k_0$  выполняется строгое неравенство

$$\sum_{k \in K} \left( \sum_{m \in M} \tau^*(m) \cdot p(L_0^n; m) \cdot p(x_0, k; m) \right) \cdot w(k, k_0) < \\ < \sum_{k \in K} \left( \sum_{m \in M} \tau^*(m) \cdot p(L_0^n; m) \cdot p(x_0, k; m) \right) \cdot w(k, k').$$

Отсюда немедленно следует, что и неравенство

$$t^*(L^{*n}, x_0) <$$

$$< \sum_{k \in K} \left( \sum_{m \in M} \tau^*(m) \cdot p(L_0^n; m) \cdot p(x_0, k; m) \right) \cdot w(k, k')$$

также выполняется строго для всех  $k' \neq k_0$ . Отсюда, в силу известной теоремы двойственности о дополняющей жесткости, следует, что  $\tilde{q}^{*n}(k_0 / L_0^n, x_0) = 1$ , т.е., что искомая стратегия есть байесовская стратегия для весовой функции  $\tau^*$ , являющейся решением двойственной задачи. Докажем теперь, что решение  $\tau^*$  двойственной задачи совпадает с весовой функцией, обеспечивающей максимальное значение числа  $\min_{q^n} R^n(q^n, m(\tau)) - \sum_{m \in M} \tau(m) \cdot R^B(m)$ , как это формулируется в теореме.

Риск  $R^n(q^n, m(\tau))$  есть по определению

$$R^n(q^n, m(\tau)) = \sum_{k' \in K} \sum_{x \in X} \sum_{L^n \in \Lambda^n} \tilde{q}^n(k' / L^n, x) \times \\ \times \sum_{k \in K} \left( \sum_{m \in M} \tau(m) \cdot p(L^n; m) \cdot p(x, k; m) \right) \cdot w(k, k').$$

Поскольку все числа  $\tilde{q}^n(k' / L^n, x)$  не отрицательные, причем для любых  $L^n \in \Lambda^n$ ,  $x \in X$  должно выполняться  $\sum_{k' \in K} \tilde{q}^n(k' / L^n, x) = 1$ , то

$$\min_{q^n} R^n(q^n, m(\tau)) = \\ = \sum_{x \in X} \sum_{L^n \in \Lambda^n} \min_{k' \in K} \sum_{k \in K} \left( \sum_{m \in M} \tau(m) \cdot p(L^n; m) \cdot p(x, k; m) \right) \times \\ \times w(k, k').$$

Поскольку в двойственной задаче величины  $t(L^n, x)$  должны принимать как можно большие значения, то их можно однозначно определить через весовые коэффициенты  $\tau(m)$ , так что

$$t(L^n, x) = \\ = \min_{k'} \sum_{k \in K} \left( \sum_{m \in M} \tau(m) p(L^n; m) \cdot p(x, k; m) \right) \cdot w(k, k')$$

и

$$\sum_{x \in X} \sum_{L^n \in \Lambda^n} t(L^n, x) - \sum_{m \in M} \tau(m) \cdot R^B(m) = \\ = \min_{q^n} R^n(q^n, m(\tau)) - \sum_{m \in M} \tau(m) R^B(m).$$

**Теорема доказана.**

Поиск оптимальной стратегии обучения-распознавания сводится, таким образом, к отысканию весовой функции  $\tau: M \rightarrow R$ , обеспечивающей максимальное значение числа

$$F(\tau) = \min_{q^n} R^n(q^n, m(\tau)) - \sum_{m \in M} \tau(m) R^B(m). \quad (20)$$

Эта задача – задача выпуклой оптимизации, как утверждает следующая теорема.

**Теорема 4.** Для любого класса  $M$  моделей функция  $F(\tau)$  является вогнутой функцией.

**Доказательство.** Докажем, что для любой весовой функции  $\tau_0 \in T$  существует линейная функция  $L_{\tau_0}: T \rightarrow R$ , такая, что неравенство

$$L_{\tau_0}(\tau - \tau_0) \geq F(\tau) - F(\tau_0) \quad (21)$$

выполняется для всех  $\tau \in T$ .

Пусть  $q_0^n$  – стратегия, обеспечивающая минимальный риск  $R^n(q^n, m(\tau_0))$ .

$$R^n(q_0^n, m(\tau_0)) = \min_{q^n} R^n(q^n, m(\tau_0)). \quad (22)$$

Определим линейную функцию  $L_{\tau_0}(\tau)$  как

$$L_{\tau_0}(\tau) = \sum_{m \in M} \tau(m) \left[ R^n(q_0^n, m) - R^B(m) \right]. \quad (23)$$

Справедлива следующая цепочка:

$$L_{\tau_0}(\tau - \tau_0) = \sum_{m \in M} \tau(m) \left[ R^n(q_0^n, m) - R^B(m) \right] - \\ - \sum_{m \in M} \tau_0(m) \left[ R^n(q_0^n, m) - R^B(m) \right] \geq \\ \geq \min_{q^n} \sum_{m \in M} \tau(m) \left[ R^n(q^n, m) - R^B(m) \right] - \\ - \sum_{m \in M} \tau_0(m) \left[ R^n(q_0^n, m) - R^B(m) \right] =$$

$$\begin{aligned}
&= \min_{q^n} \sum_{m \in M} \tau(m) [R^n(q^n, m) - R^B(m)] - \\
&- \min_{q^n} \sum_{m \in M} \tau_0(m) [R^n(q^n, m) - R^B(m)] = \\
&= F(\tau) - F(\tau_0).
\end{aligned}$$

Равенство в первом звене следует из (23). Неравенство во втором звене очевидно. Равенство в третьем звене следует из (22). Наконец, последнее равенство следует из определения (20). Таким образом, линейная функция (23) есть искомая линейная функция, доказывающая вогнутость функции  $F(\tau)$ . **Теорема доказана.**

В приведенном доказательстве получен вспомогательный результат: обобщенный градиент  $g_{\tau_0}$  функции  $F(\tau)$  в точке  $\tau_0$  есть  $|M|$ -мерный вектор с компонентами  $R^n(q_0^n, m) - R^B(m)$ . Это указывает пути отыскания весовой функции, а следовательно, и стратегии обучения методом обобщенных градиентов. В экспериментах, приводимых далее, использован другой, теоретически не обоснованный алгоритм.

Алгоритм работает по шагам  $i = 1, 2, \dots$ , и на каждом шаге строит стратегию обучения  $q_i^n$  и весовую функцию  $\tau_i = (\tau_i(m) | m \in M)$ .

Исходными для алгоритма являются конечное множество  $M$  моделей и число  $\varepsilon > 0$  – требуемая точность решения задачи.

1. Для каждой модели  $m$  определить целое число  $n(m) = 1$ .

2. Для каждой модели  $m \in M$  определить ее вес  $\tau(m) = \frac{n(m)}{\sum_{m \in M} n(m)}$ .

3. Построить байесову стратегию обучения-распознавания  $q^n$  для весов  $\tau(m)$ .

4. Для каждой модели  $m \in M$  определить число  $R^n(q^n, m)$  – качество текущей стратегии на модели  $m$ .

5. Найти гарантированное качество текущей стратегии  $R^n(q^n) = \max_m [R^n(q^n, m) - R^B(m)]$  и наихудшую модель

$$m^* = \arg \max_m [R^n(q^n, m) - R^B(m)].$$

6. Выполнить  $n(m^*) + +$ .

7. Найти среднее качество

$$\bar{R}^n(q^n) = \sum_{m \in M} \tau(m) [R^n(q^n, m) - R^B(m)].$$

8. Найти наилучший во всей предыстории алгоритма гарантированный уровень  $R^{\min} = \min(R^{\min}, R^n(q^n))$  и наилучшую во всей предыстории стратегию  $q_0^n$ .

9. Найти наихудшее во всей предыстории среднее качество  $R^{\max} = \max(R^{\max}, \bar{R}^n(q^n))$ .

10. Если  $R^{\min} - R^{\max} \leq \varepsilon$ , **останов.**

11. Перейти на п. 3.

В настоящее время теоретически не доказано, хотя экспериментально и подтверждается, что приведенный алгоритм непременно выходит на останов при любой требуемой точности  $\varepsilon > 0$ . Однако, если он вышел на останов, то его результатом является стратегия  $q_0^n$ , решающая задачу обучения с заданной точностью  $\varepsilon$  в том смысле, что

$$\max_{m \in M} [R^n(q_0^n, m) - R^B(m)] -$$

$$- \max_{m \in M} [R^n(q^{*n}, m) - R^B(m)] \leq \varepsilon,$$

где  $q^{*n}$  – решение задачи обучения (18) или, что то же самое, задачи (19).

### Иллюстрация метода обучения-распознавания

Покажем, как предложенный метод обучения-распознавания реализуется для простейшего класса моделей, описанного выше. В этом классе наблюдение  $x$  есть одномерная вещественная величина,  $K = \{1, 2\}$ , условные распределения  $p(x/k)$  полностью известны и равны

$$p(x/k) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_k)^2}, \quad k = 1, 2, \quad \mu_1 = 1, \quad \mu_2 = -1.$$

Неизвестными являются только априорные вероятности  $p_1$  и  $p_2$  состояний  $k = 1$  и  $k = 2$ . Одна конкретная модель  $m$ , таким образом, есть пара  $(p_1, p_2)$ . Поскольку в рассматриваемом примере условные распределения  $p(x/k)$  полностью известны, нет необходимости рас-

считать обучающие последовательности вида  $L^n = (x_1, k_1; x_2, k_2; \dots; x_n, k_n)$ , а достаточно рассматривать обучающие последовательности вида  $L^n = (k_1, k_2, \dots, k_n) \in \Lambda^n = K^n$ .

Для специального вида потерь  $w(k, k') = |k - k'|$ , исходя из теоремы 3, стратегия обучения-распознавания должна иметь вид

$$q\left(k' = 1 / L^n, x\right) = 1,$$

$$\text{если } \sum_{m \in M} \tau(m) \left( \prod_{i=1}^n p_{k_i} \right) \cdot p_1 \cdot p\left(x/1\right) > \\ > \sum_{m \in M} \tau(m) \left( \prod_{i=1}^n p_{k_i} \right) \cdot p_2 \cdot p\left(x/2\right);$$

$$q\left(k' = 2 / L^n, x\right) = 1,$$

$$\text{если } \sum_{m \in M} \tau(m) \left( \prod_{i=1}^n p_{k_i} \right) \cdot p_1 \cdot p(x/1) < \\ p_1 \cdot p\left(x/1\right) < \sum_{m \in M} \tau(m) \left( \prod_{i=1}^n p_{k_i} \right) \cdot p_2 \cdot p(x/2).$$

В этих выражениях для краткости не указано, что величины  $p_1$  и  $p_2$  и, вообще,  $p_{k_i}$ ,  $i = 1, 2, \dots, n$ , зависят от модели  $m$ . Указанную стратегию можно представить в более наглядном виде

$$k^* = 1, \text{ если } \frac{p(x/1)}{p(x/2)} > \theta(\tau, n_1, n_2);$$

$$k^* = 2, \text{ если } \frac{p(x/1)}{p(x/2)} < \theta(\tau, n_1, n_2),$$

где  $k^*$  – решение о состоянии  $k$  по наблюдению  $x$  и обучающей последовательности. Числа  $n_1$  и  $n_2$  обозначают, сколько раз в обучающей последовательности  $L^n$  встретилась ситуация, когда  $k_i = 1$  и  $k_i = 2$ . Решение о состоянии объекта принимается на основании сравнения отношения правдоподобия  $p(x/1)/p(x/2)$  с порогом, зависящим от обучающей выборки. Однако значение этого порога определяется как

$$\theta(\tau, n_1, n_2) = \frac{\sum_{m \in M} \tau(m) p_1^{n_1} \cdot p_2^{n_2+1}}{\sum_{m \in M} \tau(m) p_1^{n_1+1} \cdot p_2^{n_2}}, \quad (24)$$

а совсем не так просто и грубо,

$$\theta(\tau, n_1, n_2) = \frac{n_1}{n_2}, \quad (25)$$

как в методе наиболее правдоподобных моделей. По-видимому, стратегии, основанные на сравнении с порогами (24) и (25), становятся неотличимыми при неограниченном росте обучающей последовательности. Однако при малых значениях длины, как показали эксперименты, это различие существенно.

Описанный в предыдущей главе метод обучения-распознавания включает в себя и отыскание значений весов  $\tau(m)$ ,  $m \in M$ . Для реализации этого метода следует располагать конструктивным методом вычисления риска  $R^n(q^n, m)$ , который определен как

$$R^n(q^n, m) = \sum_{L^n \in \Lambda^n} \sum_{x \in X} \sum_{k' \in K} \tilde{q}^n(k' / L^n, n) p(L^n; n) \times \\ \times \sum_{k \in K} p(x, k; m) w(k, k').$$

В общем случае наибольшую трудность здесь представляет суммирование по всем обучающим выборкам  $L^n \in \Lambda^n$ . В рассматриваемом примере этот риск равен

$$R^n(q^n, m) = \sum_{L^n \in \Lambda^n} \left( \prod_{i=1}^n p_{k_i} \right) \sum_{k=1}^n p_k \cdot \Phi_k(\theta(\tau, n_1, n_2)), \quad (26)$$

где

$$\Phi_1(\theta) = \int_{-\infty}^{\theta} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-1)^2} dx, \quad \Phi_2(\theta) = \int_{\theta}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x+1)^2} dx.$$

Представим множество  $\Lambda^n$  в виде

$$\Lambda^n = \bigcup_{\substack{n_1, n_2 \\ n_1 + n_2 = n}} \Lambda(n_1, n_2),$$

где  $\Lambda(n_1, n_2)$  – множество всех тех выборок  $L^n = (k_1, k_2, \dots, k_i, \dots, k_n)$  длины  $n_1 + n_2$ , в которых событие  $k_i = 1$  произошло  $n_1$  раз, а событие  $k_i = 2 - n_2$  раз. С учетом этого обозначения выражение для риска  $R^n(q^n, m)$  приобретает вид

$$\begin{aligned}
R^n(q^n, m) &= \sum_{L^n \in \Lambda^n} p_1^{n_1} \cdot p_2^{n_2} \sum_{k=1}^2 p_k \cdot \Phi_k(\theta(\tau, n_1, n_2)) = \\
&= \sum_{\substack{n_1, n_2 \\ n_1 + n_2 = n}} \sum_{L^n \in \Lambda(n_1, n_2)} p_1^{n_1} \cdot p_2^{n_2} \sum_{k=1}^2 p_k \cdot \Phi_k(\theta(\tau, n_1, n_2)) = \\
&= \sum_{\substack{n_1, n_2 \\ n_1 + n_2 = n}} p_1^{n_1} \cdot p_2^{n_2} \sum_{k=1}^2 p_k \cdot \Phi_k(\theta(\tau, n_1, n_2)) \sum_{L^n \in \Lambda(n_1, n_2)} 1 = \\
&= \sum_{\substack{n_1, n_2 \\ n_1 + n_2 = n}} C_n^{n_1} \cdot p_1^{n_1+1} \cdot p_2^{n_2} \Phi_1(\theta(\tau, n_1, n_2)) + \\
&+ \sum_{\substack{n_1, n_2 \\ n_1 + n_2 = n}} C_n^{n_1} \cdot p_1^{n_1} \cdot p_2^{n_2+1} \Phi_2(\theta(\tau, n_1, n_2)),
\end{aligned}$$

и его вычисление перестает быть проблематичным.

На рис. 2–7 показана зависимость риска различных стратегий от модели  $m$ , в данном случае, от вероятности  $p_1$  первого состояния. Байесовский риск  $R^B(m)$  показан сплошной линией, риск стратегии по методу наибольшего правдоподобия показан пунктирной линией, а риск стратегии по предлагаемому методу – жирной пунктирной линией. Рис. 2–7 соответствуют различным значениям длины  $n$  обучающих выборок: 1, 2, 3, 10, 20 и 50.

Расчет весов  $\tau(m)$ ,  $m \in M$ , выполнялся в соответствии с алгоритмом, описанным в предыдущем разделе.

**Заключение.** Результаты данной статьи изложены на определенном уровне общности,

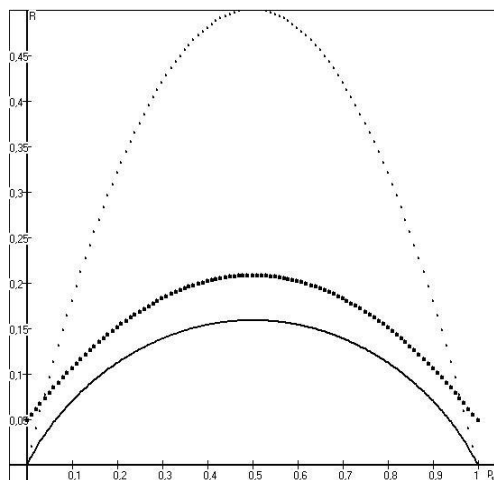


Рис. 2. Вероятности ошибочного распознавания для различных стратегий при длине обучающей выборки, равной 1

который мог быть и большим, и меньшим. Обозначим некоторые возможные обобщения приведенных результатов, которые уместно рассмотреть в данной статье.

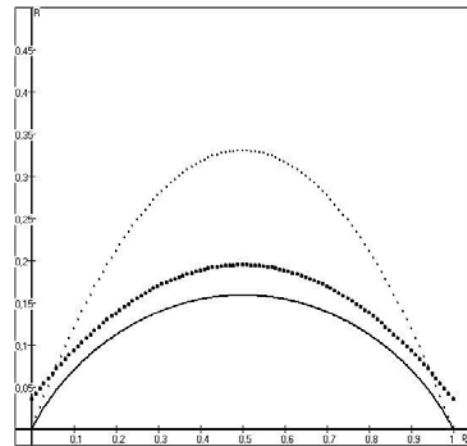


Рис. 3. Вероятности ошибочного распознавания для различных стратегий при длине обучающей выборки, равной 2

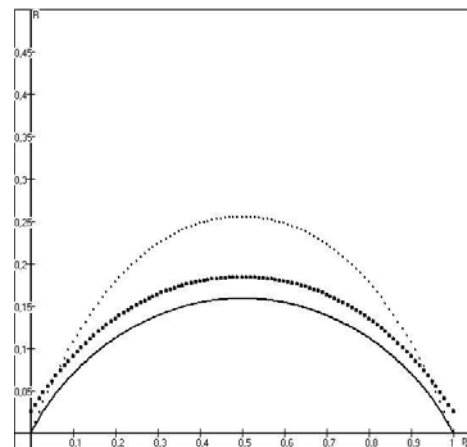


Рис. 4. Вероятности ошибочного распознавания для различных стратегий при длине обучающей выборки, равной 3

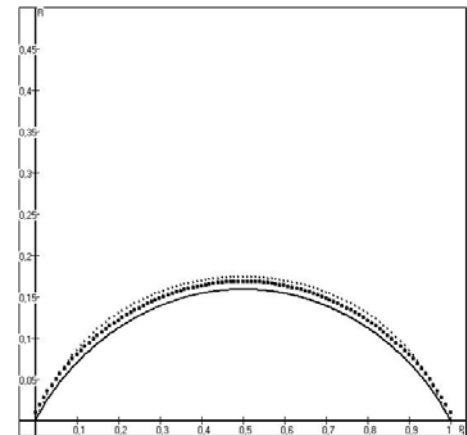


Рис. 5. Вероятности ошибочного распознавания для различных стратегий при длине обучающей выборки, равной 10

## Обучение и самообучение

Обучающая информация, на основании которой строится собственно распознавание, совсем не обязательно должна быть выборкой вида  $(x_1, k_1; x_2, k_2; \dots; x_n, k_n) \in (X \times K)^n$ . Это может быть результат  $L$  любого эксперимента над рас-

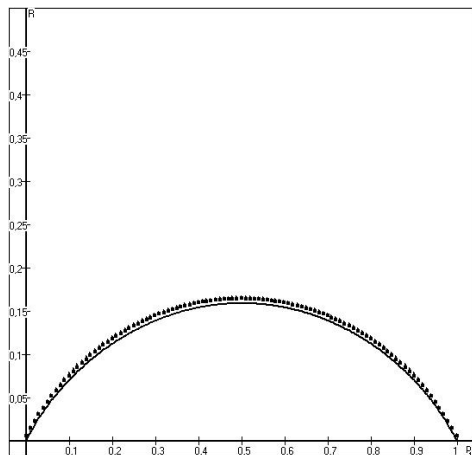


Рис. 6. Вероятности ошибочного распознавания для различных стратегий при длине обучающей выборки, равной 20

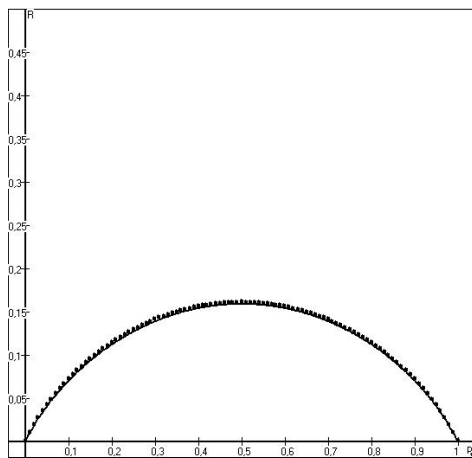


Рис. 7. Вероятности ошибочного распознавания для различных стратегий при длине обучающей выборки, равной 50

познаваемым объектом. Важно только, чтобы этот результат мог считаться случайным и зависящим от модели, так, чтобы имела смысл вероятность  $p(L; m)$ . Результаты статьи без изменения переносятся и на этот более общий случай. В частности, обучающая информация может быть выборкой  $L = (x_1, x_2, \dots, x_n) \in X^n$ , состоящей только из наблюдаемых при распознавании признаков и не включающей в себя

информацию о скрытых состояниях. Построение распознающих систем при обучающей информации такого вида получило название *самообучение*. Основным методом построения таких систем стало, как и при обучении, максимально правдоподобное оценивание модели на основании выборки  $L = (x_1, x_2, \dots, x_n)$ , для которого широко применяются алгоритмы [5], ставшие известными как *EM*-алгоритмы [6].

Описанный подход вносит коррективы и в эти известные методы распознавания в режиме самообучения. Как и в случае обучения, любые разумные требования к распознаванию в режиме самообучения приводят к алгоритмам, в которых отсутствует такой этап, как поиск наиболее правдоподобной модели на основании выборки  $L = (x_1, x_2, \dots, x_n)$ . Эти коррективы могут быть существенными, особенно при малых значениях длины  $n$  выборки  $L$ .

Здесь возможны и более существенные выводы, которые покажем на примере. Пусть для заданного класса моделей и обучающей информации произвольного вида построена стратегия  $\tilde{q}(k/L, x)$  обучения-распознавания. Коль скоро эта стратегия построена, то она предназначена не для однократного применения при распознавании какого-то одного объекта, а для распознавания многих объектов. При этом предполагается, что все эти объекты характеризуются той же моделью  $m$ , при которой была получена обучающая информация  $L$ . Пусть теперь требуется принять решение о состояниях  $k_1$  и  $k_2$  двух объектов по признакам  $x_1$  и  $x_2$  этих объектов. В этом случае можно поступить двумя способами.

*Первый* способ состоит в применении полученной стратегии  $\tilde{q}(k/L, x)$  к каждому из двух предъявленных объектов и с вероятностью  $\tilde{q}(k_1/L, x_1)$  принять решение, что первый объект находится в состоянии  $k_1$ , и с вероятностью  $\tilde{q}(k_2/L, x_2)$  принять решение, что второй объект находится в состоянии  $k_2$ .

*Второй* способ состоит в рассмотрении пары предъявленных объектов как одного объекта с

состоянием из множества  $K \times K = \{(k_1, k_2) \mid k_1 \in K, k_2 \in K\}$ , с признаком  $(x_1, x_2)$ , который принимает значения из множества  $X \times X$ , и классом  $M$  моделей  $m$ , таких, что

$$p(x_1, x_2, k_1, k_2; m) = p(x_1, k_1; m) \cdot p(x_2, k_2; m). \quad (27)$$

Для такого класса моделей строится стратегия  $q'(k_1, k_2 / L, x_1, x_2)$ , в соответствии с которой решения о состояниях  $(k_1, k_2)$  пары объектов принимаются на основании пары  $(x_1, x_2)$  признаков. Эти решения совсем не обязательно окажутся независимыми одно от другого. Из предположения (27) не следует, что

$$q'(k_1, k_2 / L, x_1, x_2) = q(k_1 / L, x_1) \cdot q(k_2 / L, x_2).$$

Даже в случае, когда два объекта, предъявленные для распознавания, взаимно независимы, решения о их состояниях не являются независимыми. При этом качество стратегии  $q'$ , конечно, не хуже, чем качество стратегии  $q$ , а может быть и значительно лучшим.

Таким образом обнаруживается еще одна предпосылка, принятая в современных подходах к обучению, как сама собою разумеющаяся, но не вытекающая из каких-либо разумных соображений. Считается, что на основании обучающей выборки происходит выбор той или иной стратегии собственно распознавания, которая затем в неизменном виде применяется для всех объектов, предъявленных для распознавания. Это – ненужное и ниоткуда не следующее ограничение.

Необходимость в обучении возникает вообще только при априорной неопределенности модели объекта. Обучающая информация позволяет лишь уменьшить эту неопределенность, а не исключить ее. Каждое из наблюдений, предъявленных для распознавания уже после обработки обучающей информации, статистически зависит от модели и поэтому дополнительно уменьшает ее неопределенность. Если для распознавания предъявлена последовательность из  $l$  объектов, т.е. последовательность  $x_1, x_2, \dots, x_l$  признаков этих объектов, то только первый объект должен распознаваться в соответствии со стратегией  $q(k_1 / L, x_1)$ . При распознавании второго объекта признак  $x_1$  стано-

вится уже частью обучающей информации и распознавание должно выполняться в соответствии с другой стратегией вида  $q((k_2 / (L, x_1), x_2)$ . И вообще,  $i$ -й объект должен распознаваться в соответствии со стратегией вида  $q(k_i / (L, x_1, \dots, x_{i-1}), x_i)$ .

Обучение, т.е. изменение стратегии распознавания, совсем не должно заканчиваться обработкой обучающей информации  $L$ . Обработка этой информации является лишь нулевым этапом обучения, которое затем продолжается в течение всего периода эксплуатации распознающей системы.

### **Эмпирический байесовский подход**

Очевидно, что при неполном знании статистической модели распознавание многих объектов одновременно может оказаться заметно лучшим, чем независимое одно от другого распознавание каждого объекта в отдельности. Этот вывод справедлив независимо от вида обучающей последовательности, в том числе и в ее отсутствие. Эту идею в частном виде высказал Г. Роббинс при формулировке предложенного им эмпирического байесовского подхода [7, 8]. Ученый сформулировал его для случая, когда модель распознаваемого объекта почти полностью известна, а неизвестны только априорные вероятности скрытых состояний. Для этого примера была указана реализация предлагаемого подхода при неограниченном росте количества объектов, предъявленных для распознавания. Асимптотический характер подхода выражен в самом его названии «асимптотически субминимаксные решения в составных задачах теории статистических решений». Вопрос о том, каким должно быть решение не в асимптотическом случае, а при распознавании конечных совокупностей объектов, был оставлен для будущих исследований и до сих пор оставался без ответа. Исследования задачи обучения указывают путь к реализации эмпирического байесовского подхода и в случае конечной совокупности объектов, предъявленных для распознавания.

### **Проблема коротких выборок**

В изложенном подходе совсем не обязательно считать, что распознавание объекта состоит

в указании непременно состояния объекта. Можно допустить и специальный ответ распознающей системы, обозначаемый символом  $\# \notin K$  и понимаемый как отказ от распознавания или ответ НЕ ЗНАЮ. В этом случае функция потерь имеет вид  $w: K \times (K \cup \{\#\}) \rightarrow R$ , а не  $w: K \times K \rightarrow R$ , как раньше, а стратегия обучения-распознавания принимает вид  $\tilde{q}^n: L^n \times X \times (K \cup \{\#\}) \rightarrow R$ , а не  $\tilde{q}^n: L^n \times X \times K \rightarrow R$ .

Изложенные идеи построения стратегии обучения-распознавания без каких-либо затруднений можно обоснованно изложить и при таком расширенном определении основных понятий.

Введение в рассмотрение понятия «отказ от распознавания» приводит к новой точке зрения на проблему коротких выборок. Исследования по проблеме таких выборок имеют целью ответить на вопрос, при какой минимальной длине обучающей выборки последующее распознавание остается еще удовлетворительным в том или ином смысле. Все выборки с длиной, ниже этого предельного значения, объявляются короткими и непригодными для последующего распознавания. Фактически это означает отказ от распознавания независимо от того, какой именно оказалась конкретная короткая выборка и каким оказалось текущее наблюдение над объектом, предъявленным для распознавания.

Формулировка задачи обучения как построение оптимальной стратегии обучения-распознавания, допускающей и отказ от распознавания, позволяет принимать решение о хороших или плохих выборках не просто на основании значений их длины, а более точно, на основании того, какой именно оказалась та или иная длинная или короткая выборка. А именно, среди выборок, пусть даже очень коротких, есть такие, на основании которых возможно вполне удовлетворительное последующее распознавание, равно как среди длинных могут оказаться выборки, неудачные в этом отношении.

Понимание того, что результатом обучения есть стратегия обучения-распознавания, по-

зволяет идти еще дальше. А именно, никакая выборка, пусть даже очень короткая или очень длинная, сама по себе не является ни хорошей, ни плохой. Любая выборка позволяет качественно распознавать одни наблюдения и не годится для распознавания других. Хорошей или плохой является не выборка сама по себе, а пара «выборка-наблюдение», и стратегия обучения-распознавания через решение «отказ от распознавания» реализует в том числе и эту классификацию.

Таким образом, классификация выборок на короткие, т.е. плохие, и длинные, т.е. хорошие, есть лишь сильно загрубленная постановка вопроса о построении оптимальной стратегии обучения-распознавания с возможностью отказа от распознавания. Построение оптимальной стратегии обучения-распознавания поглощает проблему коротких выборок и, поглотив ее, решает ее более тонко, более дифференцированно.

1. Шлезингер М., Главач В. Десять лекций по статистическому и структурному распознаванию. – К.: Наук. думка, 2004. – 545 с.
2. Гупал А.М., Пашко С.В., Сергиенко И.В. Эффективность байесовской процедуры распознавания // Кибернетика и системный анализ. – 1995. – № 4. – С. 76–89.
3. Гупал А.М., Сергиенко И.В. Оптимальные процедуры распознавания и их применение // Там же. – 2007. – № 6. – С. 41–54.
4. Зуховицкий С.И., Авдеева Л.И. Линейное и выпуклое программирование. – М.: Наука, 1967. – 460 с.
5. Шлезингер М.И. Взаимосвязь обучения и самообучения в распознавании образов // Кибернетика. – 1968. – № 2. – С. 81–88.
6. Demster A., Laird N., Rubin D. Maximum likelihood from incomplete data via the EM algorithm // J. of the royal Statistic Society. – 1977. – В39. – Р. 1–38.
7. Роббинс Г. Асимптотически субминимаксные решения в составных задачах теории статистических решений // Математика. – 1964. – № 8:2. – С. 141–159.
8. Роббинс Г. Эмпирический байесовский подход к статистике // Там же. – С. 133–140.