

## Задача обучения распознаванию образов в нестационарной генеральной совокупности

*Шавловский М. Б., Красоткина О. В., Моттль В. В.*

shavlovsky@yandex.ru, ko180177@yandex.ru, vmottl@yandex.ru

Москва, МФТИ; Тула, ТулГУ; Москва, ВЦ РАН

Целью данного исследования является создание основного математического аппарата и простейших алгоритмов для решения типичных для практики задач обучения распознаванию образов в генеральных совокупностях, свойства которых изменяются во времени. Широко известная классическая постановка задачи распознавания основана на молчаливом предположении, что свойства генеральной совокупности в момент «экзамена» остаются теми же, что и при формировании обучающей выборки. Принятое в данной работе более реалистичное предположение о нестационарности генеральной совокупности неизбежно приводит к необходимости анализа последовательности выборок в некоторые моменты времени и поиска для них разных решающих правил распознавания.

Пусть каждый объект генеральной совокупности  $\omega \in \Omega$  представлен точкой в линейном пространстве признаков  $\mathbf{x}(\omega) = (x^1(\omega), \dots, x^n(\omega)) \in \mathbb{R}^n$ , а его скрытая фактическая принадлежность к одному из двух классов определяется значением индекса класса  $y(\omega) \in \{1, -1\}$ . Классический подход к обучению распознаванию двух классов объектов, развитый В. Н. Вапником [1], основан на понимании модели генеральной совокупности в виде дискриминантной функции  $f(\mathbf{x}(\omega)) = \mathbf{a}^T \mathbf{x} + b$ , определяемой гиперплоскостью с априори неизвестными наблюдателю направляющим вектором  $\mathbf{a} \in \mathbb{R}^n$  и параметром положения  $b \in \mathbb{R}$ :

$$f(\mathbf{x}(\omega)) = \mathbf{a}^T \mathbf{x} + b \begin{cases} \text{преимущественно} > 0, & \text{если } y(\omega) = 1, \\ \text{преимущественно} < 0, & \text{если } y(\omega) = -1. \end{cases}$$

Неизвестные параметры разделяющей гиперплоскости подлежат оцениванию на основе анализа обучающей совокупности объектов  $\{\omega_j, j = 1, \dots, N\}$ , представленных векторами их признаков и индексами принадлежности к классам, так что выборка в целом является конечным множеством пар  $\{(\mathbf{x}_j \in \mathbb{R}^n, y_j = \pm 1), j = 1, \dots, N\}$ . Широко известен принцип оптимальной разделяющей гиперплоскости, выбираемой по критерию максимизации числа точек обучающей совокупности, правильно классифицируемых с гарантированным «запасом», условно равным единице:

$$\begin{cases} J(\mathbf{a}, b, \delta_j) = \mathbf{a}^T \mathbf{a} + C \sum_{j=1}^N \delta_j \rightarrow \min, \\ y_j(\mathbf{a}^T \mathbf{x}_j + b) \geq 1 - \delta_j, \quad \delta_j \geq 0, \quad j = 1, \dots, N. \end{cases} \quad (1)$$

Понятие времени здесь полностью отсутствует.

Принципиальное отличие предлагаемой в данной работе концепции нестационарной генеральной совокупности заключается во введении в рассмотрение фактора времени  $t$ . Предполагается, что основное свойство нестационарной генеральной совокупности выражается изменяющейся во времени разделяющей гиперплоскостью, характеризующей преимущественное различие векторов признаков объектов двух классов и, в свою очередь, полностью определяемой своим направляющим вектором и параметром положения, которые должны рассматриваться как функции времени  $\mathbf{a}_t$  и  $b_t$ :

$$f_t(\mathbf{x}(\omega)) = \mathbf{a}_t^T \mathbf{x} + b_t \begin{cases} \text{преимущественно } > 0 \text{ в момент } t, & \text{если } y(\omega) = 1, \\ \text{преимущественно } < 0 \text{ в момент } t, & \text{если } y(\omega) = -1. \end{cases}$$

Здесь всякий объект  $\omega \in \Omega$  рассматривается всегда только вместе с указанием момента времени, в который он предъявлен  $(\omega, t)$ . В результате обучающая совокупность приобретает структуру множества троек  $\{(\mathbf{x}_j \in \mathbb{R}^n, y_j = \pm 1, t_j), j = 1, \dots, N\}$ , а не пар. Естественно нумеровать объекты обучающей совокупности в порядке поступления объектов, тогда уместно говорить скорее об обучающей последовательности, нежели об обучающей совокупности, рассматривая ее как временной ряд, вообще говоря, с переменным шагом по времени.

В разный момент времени  $t_j$  скрытая от наблюдателя разделяющая гиперплоскость характеризуется разными неизвестными значениями направляющего вектора и параметра положения. Таким образом, объективно существует двухкомпонентный временной ряд со скрытой и наблюдаемой компонентами, соответственно,  $(\mathbf{a}_j, b_j)$  и  $(\mathbf{x}_j, y_j)$ .

В динамической постановке задача обучения превращается в задачу анализа двухкомпонентного временного ряда, в котором требуется, анализируя наблюдаемую компоненту, дать оценку скрытой компоненты. Это стандартная задача анализа сигналов (временных рядов), специфика которой заключается лишь в предполагаемой модели связи между скрытой и наблюдаемой компонентами. Согласно классификации задач оценивания скрытой компоненты сигнала, введенной Н. Винером [2], естественно различать, по крайней мере, два вида задач обучения.

*Задача фильтрации обучающей последовательности.* Пусть  $t_j$  — момент поступления очередного объекта, к которому уже зарегистрированы векторы признаков и индексы классов объектов  $\{\dots, (\mathbf{x}_{j-2}, y_{j-2}), (\mathbf{x}_{j-1}, y_{j-1}), (\mathbf{x}_j, y_j)\}$ , поступивших в предыдущие моменты времени до текущего момента включительно  $(\dots, t_{j-2}, t_{j-1}, t_j)$ . Требуется непосредственно в процессе наблюдения давать оценку параметров разделяющей гиперплоскости  $(\hat{\mathbf{a}}_j, \hat{b}_j)$  в каждый текущий момент времени  $t_j$ .

*Задача интерполяции.* Пусть к моменту обработки обучающая последовательность уже зарегистрирована в некотором интервале времени  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ . Требуется оценить изменяющиеся параметры разделяющей гиперплоскости во всем интервале наблюдения  $\{(\hat{\mathbf{a}}_1, \hat{b}_1), \dots, (\hat{\mathbf{a}}_N, \hat{b}_N)\}$ .

Предполагается, что параметры разделяющей гиперплоскости  $\mathbf{a}_t$  и  $b_t$  изменяются во времени достаточно медленно в том смысле, что величины

$$\frac{1}{t_j - t_{j-1}}(\mathbf{a}_j - \mathbf{a}_{j-1})^T(\mathbf{a}_j - \mathbf{a}_{j-1}) \cong \varepsilon_a \quad \text{и} \quad \frac{1}{t_j - t_{j-1}}(b_j - b_{j-1})^2 \cong \varepsilon_b$$

являются, как правило, достаточно малыми. Это предположение препятствует вырождению задач фильтрации и интерполяции в совокупность независимых некорректных задач обучения распознаванию двух классов объектов по единственному наблюдению.

С формальной точки зрения оценка параметров разделяющей гиперплоскости в последний момент интервала наблюдения  $(\hat{\mathbf{a}}_N, \hat{b}_N)$ , полученная при решении задачи интерполяции, является решением задачи фильтрации для этого момента. Однако смысл задачи фильтрации заключается в том, чтобы очередные оценки вычислялись непосредственно в процессе поступления новых наблюдений, без решения всякий раз задачи интерполяции для временного ряда возрастающей длины.

Предлагаемая постановка задачи обучения в режиме интерполяции отличается от совокупности классических задач обучения по методу опорных векторов (1) для каждого момента времени только наличием дополнительных членов, штрафующих различие между смежными значениями параметров гиперплоскости  $(\mathbf{a}_{j-1}, b_{j-1})$  и  $(\mathbf{a}_j, b_j)$ :

$$\begin{cases} J(\mathbf{a}_j, b_j, \delta_j, j = 1, \dots, N) = \sum_{j=1}^N \left( \frac{1}{N} \mathbf{a}_j^T \mathbf{a}_j + \delta_j \right) + \\ + \sum_{j=2}^N \frac{1}{t_j - t_{j-1}} [D^a (\mathbf{a}_j - \mathbf{a}_{j-1})^T (\mathbf{a}_j - \mathbf{a}_{j-1}) + D^b (b_j - b_{j-1})^2] \rightarrow \min, \\ y_j (\mathbf{a}_j^T \mathbf{x}_j + b_j) \geq 1 - \delta_j, \quad \delta_j \geq 0, \quad j = 1, \dots, N. \end{cases} \quad (2)$$

Здесь коэффициенты  $D^a > 0$  и  $D^b > 0$  являются параметрами критерия, задающими желаемую степень сглаживания оцениваемой последовательности мгновенных значений параметров разделяющей гиперплоскости.

Критерий (2) реализует концепцию оптимальной достаточно гладкой последовательности разделяющих гиперплоскостей в отличие от концепции единственной оптимальной гиперплоскости в (1). Искомые гиперплоскости должны обеспечивать правильную классификацию векторов признаков объектов для как можно большего числа моментов времени с гарантированным «запасом», принятым равным единице, как и в (1).

Как и классическая задача обучения, динамическая задача (2) является задачей квадратичного программирования, но содержит  $N(n+1)+N$  переменных, в отличие от  $(n+1)+N$  переменных в (1). Известно, что вычислительная сложность задачи квадратичного программирования общего вида пропорциональна кубу числа переменных, т. е. динамическая задача, на первый взгляд, существенно сложнее классической.

Однако целевая функция  $J(\mathbf{a}_j, b_j, \delta_j, j = 1, \dots, N)$  в динамической задаче является парно-сепарабельной, т. е. представляя собой сумму частных функций, каждая из которых зависит от переменных, связанных только с одним либо двумя моментами времени в порядке их возрастания. Это обстоятельство позволяет построить алгоритм численного решения задачи, вычислительная сложность которого линейна относительно длины обучающей последовательности  $N$ .

Применение теоремы Куна-Таккера к динамической задаче (2) переводит ее в двойственную форму относительно множителей Лагранжа  $\lambda_j \geq 0$  при ограничениях-неравенствах  $y_j(\mathbf{a}_j^T \mathbf{x}_j + b_j) \geq 1 - \delta_j$ :

$$\begin{cases} W(\lambda_1, \dots, \lambda_N) = \\ = \sum_{j=1}^N \lambda_j - \frac{1}{2} \sum_{j=1}^N \sum_{l=1}^N y_j y_l (\mathbf{a}_j^T \mathbf{Q}_{jl} \mathbf{a}_l + f_{jl}) \lambda_j \lambda_l \rightarrow \max, \\ \sum_{j=1}^N y_j \lambda_j = 0, \quad 0 \leq \lambda_j \leq C/2, \quad j = 1, \dots, N. \end{cases} \quad (3)$$

Здесь матрицы  $\mathbf{Q}_{jl}$  ( $n \times n$ ) и  $\mathbf{F} = (f_{jl})$  ( $N \times N$ ) не зависят от обучающей последовательности и определяются только коэффициентами штрафа, соответственно,  $D^a$  на негладкость последовательности направляющих векторов искомого гиперплоскостей и  $D^b$  на негладкость последовательности их параметров положения в (2).

**Теорема 1.** *Решение задачи обучения (2) полностью определяется значениями множителей Лагранжа  $(\lambda_1, \dots, \lambda_N)$ , полученными как решение двойственной задачи (3), и обучающей последовательностью:*

$$\hat{\mathbf{a}}_j = \sum_{l: \lambda_l > 0} y_l \lambda_l \mathbf{Q}_{jl} \mathbf{x}_l; \quad \hat{b}_j = b + \sum_{l: \lambda_l > 0} y_l \lambda_l f_{jl}; \quad (4)$$

$$b = \frac{\sum_{j: 0 < \lambda_j < C/2} \lambda_j \sum_{l: \lambda_l > 0} y_l \lambda_l (\mathbf{x}_j^T \mathbf{Q}_{jl} \mathbf{x}_j + f_{jl}) + (C/2) \sum_{j: \lambda_j = C/2} y_j}{\sum_{j: 0 < \lambda_j < C/2} \lambda_j}. \quad (5)$$

Из этих формул видно, что решение задачи динамического обучения определяется только теми элементами обучающей последовательности

$(\mathbf{x}_j, y_j)$ , множители Лагранжа при которых получили положительные значения  $\lambda_j > 0$ . Уместно назвать векторы признаков соответствующих объектов опорными векторами, так что мы пришли к некоторому обобщению метода опорных векторов [1], вытекающего из концепции оптимальной разделяющей гиперплоскости (1).

Классическая задача обучения (1) является частным случаем задачи (2) при бесконечно больших значениях штрафов на изменение параметров гиперплоскости  $D^a \rightarrow \infty$  и  $D^b \rightarrow \infty$ . В этом случае  $\mathbf{Q}_{jl} \rightarrow \mathbf{I}$ ,  $f_{jl} \rightarrow 0$ , и двойственная задача (3) превращается в классическую двойственную задачу [1], соответствующую исходной задаче (1), а формулы (4) и (5) определяют результат обучения по классическому методу опорных векторов  $\hat{\mathbf{a}} = \hat{\mathbf{a}}_1 = \dots = \hat{\mathbf{a}}_N$ ,  $\hat{b} = \hat{b}_1 = \dots = \hat{b}_N$ .

Хотя двойственная задача (3) и не является парно-сепарабельной, в силу парной сепарабельности исходной задачи (2) вычисление градиента целевой функции  $\nabla_{\boldsymbol{\lambda}} W(\lambda_1, \dots, \lambda_N)$  в любой точке  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_N)$  и, далее, определение оптимального допустимого направления поиска с учетом ограничений обеспечиваются алгоритмом, имеющим линейную вычислительную сложность относительно длины обучающей последовательности. В результате оказывается, что использование любого градиентного метода для решения двойственной задачи приводит к алгоритму с линейной вычислительной сложностью относительно  $N$ . В частности, стандартный метод наискорейшего спуска для решения задач квадратичного программирования [3], примененный к функции  $-W(\lambda_1, \dots, \lambda_N)$ , дает алгоритм, являющийся, по сути, обобщением известного алгоритма SMO (Sequential Minimum Optimization) [4], обычно используемого при решении двойственных задач.

Работа выполнена при поддержке РФФИ, проекты № 05-01-00679, № 06-01-00412.

### Литература

- [1] *Vapnik V.* Statistical Learning Theory. — New York: John-Wiley & Sons, Inc., 1998. — 732 p.
- [2] *Wiener N.* Extrapolation, Interpolation, and Smoothing of Stationary Random Time Series with Engineering Applications. — Technology Press of MIT, John Wiley & Sons, 1949. — 163 p.
- [3] *Базара М., Шетти К.* Нелинейное программирование. Теория и алгоритмы. — М.: Мир, 1982.
- [4] *Platt J.C.* Fast training of support vector machines using sequential minimal optimization / Advances in Kernel Methods: Support Vector Learning. — MIT Press, Cambridge, MA, 1999.