

УДК 004.93

С.А. Зайцев¹, С.А. Субботин²¹ Запорожский национальный технический университет, г. Запорожье, zaitsev.serge@gmail.com² Запорожский национальный технический университет, г. Запорожье, subbotin@zntu.edu.ua

МОДЕЛЬ ОТРИЦАТЕЛЬНОГО ОТБОРА С ИСПОЛЬЗОВАНИЕМ МАСКИРОВАННЫХ ДЕТЕКТОРОВ И МЕТОД ЕЁ ОБУЧЕНИЯ ДЛЯ РЕШЕНИЯ ЗАДАЧ ДИАГНОСТИРОВАНИЯ

Исследовалось применение модели отрицательного отбора в диагностировании. Предложена модель отрицательного отбора, использующая маскирование детекторов, и разработан метод ее обучения, что позволило повысить скорость работы модели и улучшить интерпретируемость получаемых результатов. Проведены эксперименты, подтверждающие эффективность предложенной модели.

ОТРИЦАТЕЛЬНЫЙ ОТБОР, БИТОВАЯ СТРОКА, ДЕТЕКТОР, ПРАВИЛО СОПОСТАВЛЕНИЯ, МАСКИРОВАНИЕ

Введение

При решении задач технического и медицинского диагностирования возникает необходимость определения дефектных изделий или опасных состояний объектов диагностирования, что предполагает наличие диагностической модели. Для обеспечения удобства применения диагностическую модель целесообразно представлять в виде набора продукционных правил вида “если-то”, которые могут быть получены посредством индуктивного обучения по прецедентам.

В простейшем случае каждое правило основывается на бинарном представлении антецедентов, где значением “1” кодируется наличие признака, а значением “0” – его отсутствие у данного экземпляра. Исходя из значений битовой строки, описывающей экземпляр, принимается решение об отнесении его к классу годных или дефектных.

Используемые для решения данной задачи подходы имеют ряд недостатков. В частности, нейронные сети требуют дополнительной процедуры вербализации для упрощения правил [1]. Деревья принятия решений часто сходятся на локальных оптимумах, при большом количестве признаков полученные правила значительно усложняются, а также деревья решений плохо поддаются переобучению [2].

В целях устранения перечисленных выше недостатков при решении данной задачи целесообразно использовать принцип отрицательного отбора в искусственных иммунных системах. Он обладает таким рядом свойств, как способность работать с бинарным представлением признаков, возможность обучаться на экземплярах только одного класса, распределенность вычислений [3].

Среди недостатков существующих реализаций модели отрицательного отбора [4] стоит отметить тот факт, что они требуют предварительной оценки и отбора информативных признаков, а также характеризуются высокой сложностью извлечения знаний из полученных результатов работы модели.

Цель работы заключается в разработке такой модели отрицательного отбора, которая бы позволяла проводить отбор информативных групп признаков и формировать на их основе продукционные правила для проведения диагностирования, а также разработать метод ее обучения.

1. Постановка задачи

Пусть мы имеем выборку S' , состоящую из экземпляров, описанных битовыми строками фиксированной длины l . Будем считать, что набор всех возможных битовых строк длиной l формирует пространство признаков U . Множество U можно разделить на два комплементарных подмножества, описывающих “свои” (годные) и “чужие” (дефектные) экземпляры соответственно: $U = S \cup N$, $S \cap N = \emptyset$, где S – множество годных экземпляров, а N – множество дефектных. Обучающая выборка состоит только из годных экземпляров $S' \subset S$. В этом случае задача построения модели отрицательного отбора заключается в генерации такого набора правил, представленного множеством детекторов D , на основании которого каждый $x \in U$ можно однозначно отнести к классу годных или дефектных экземпляров.

2. Метод обучения модели отрицательного отбора с цензурированием

Рассмотрим модель [5], реализующую парадигму отрицательного отбора.

Как правило, для определения принадлежности экземпляра к множеству S или N модель отрицательного отбора в процессе обучения добавляет в набор D такие детекторы, которые не соответствуют “своим” экземплярам. Поскольку множества S и N комплементарны, то предполагается, что любая битовая строка x принадлежит множеству N , если ей соответствует хотя бы один детектор из набора D . Определение соответствия экземпляра детектору происходит на основании правила сопоставления $match(d, x)$, которое принимает значение “истина”, если детектор соответствует экземпляру, и “ложь” – в противном случае.

Таким образом, работа данной модели осуществляется в два этапа.

1. Генерация набора детекторов. Для этого случайно сгенерированные кандидаты в детекторы C , представленные в виде битовых строк, подлежат цензурированию, и те из них, которые не отсеиваются на данном этапе, попадают в набор детекторов D . В результате цензурирования отсеиваются те кандидаты в детекторы $c \in C$, для которых существует такой $x \in S$, чтобы $match(c, x) = 1$.

2. Классификация. На этом этапе экземпляр x , поступающий на вход модели, сравнивается с детекторами из набора D , используя правило сопоставления. Если хотя бы один из детекторов при этом активизируется, т.е. $\exists d \in D: match(d, x) = 1$, считается, что $x \in N$, в противном случае — $x \in S$.

На практике оказывается [5], что множество D относительно небольшой мощности способно обеспечить достаточно высокую точность классификации диагностируемых данных. Более того, при неизменном количестве детекторов объем годных экземпляров может увеличиваться без снижения точности диагностирования.

Метод генерации детекторов представляется чрезвычайно ресурсоемким, а потому очень важно своевременно осуществить останов во избежание генерации избыточного количества детекторов.

3. Бинарные метрики

В качестве правила сопоставления двух бинарных детекторов применяются различного рода метрики, позволяющие определить степень подобия двух битовых строк [5-8].

Правило г-последовательных битов (*g-contiguous rule*, *RCB rule*) [5-7] использовалось изначально в модели отрицательного отбора. Метод генерации детекторов для модели отрицательного отбора оперировал строками фиксированной длины. Так, для двух строк, представленных в виде последовательности из n битов $x = \{x_1, x_2, \dots, x_n\}$ и $d = \{d_1, d_2, \dots, d_m\}$ правило выглядит следующим образом:

$$match(d, x) = \begin{cases} 1, \exists i, i \leq n - r + 1, \forall i \leq j \leq i + r - 1: x_j = d_j; \\ 0, \text{ в противном случае.} \end{cases}$$

Иными словами, две строки совпадают, если существует такое окно размером r , в пределах которого все биты обеих строк совпадают.

Данная метрика отличается своей простотой и используется в оригинальном методе генерации детекторов для модели отрицательного отбора, получившем свое дальнейшее развитие в линейном и “жадном” методах генерации детекторов [6]. Все эти методы ограничиваются использованием бинарной формы представления детекторов и RCB-правила.

Метрика R-chunks [7] является более общей по сравнению с RCB-метрикой, т.е. любой детектор,

оперирующий RCB-правилом, может быть представлен в виде множества r -chunks детекторов. Для двух строк $x = \{x_1, x_2, \dots, x_n\}$ и $d = \{d_1, d_2, \dots, d_m\}$ длиной n и m соответственно, $n \leq m$, правило r -chunks можно представить как:

$$match(d, x) = \begin{cases} 1, \forall i \leq j \leq i + m - 1: x_j = d_j; \\ 0, \text{ в противном случае,} \end{cases}$$

где i определяет позицию начала отрезка строки (chunk).

Правило r -chunks позволяет повысить точность работы метода отрицательного отбора.

Метрика Хемминга применялась в [8] в качестве правила сопоставления:

$$match(d, x) = \begin{cases} 1, \sum_{i=1}^n |x_i \oplus d_i| \geq r; \\ 0, \text{ в противном случае,} \end{cases}$$

где n — длина битовой строки, \oplus — операция “исключающее ИЛИ”, r — порог срабатывания правила, $0 \leq r \leq n$.

Метрика Левенштейна [8] может считаться обобщением метрики Хемминга. Её значение определяется минимальным количеством изменений, необходимых для преобразования одной бинарной строки в другую. При этом изменения могут быть следующего вида: вставка разряда, удаление разряда, замена одного разряда (бита) другим из алфавита. В некоторых вариациях метрика Левенштейна считает замену нескольких смежных разрядов одной операцией.

Использование метрики Левенштейна целесообразно, если экземпляры обладают различным количеством признаков.

Важно отметить, что в случае использования метрики *gcb* и *g-chunks* необходимо заранее учитывать информативность признаков и возможно произвести перестановку битов в строках таким образом, чтобы биты, соответствующие информативным признакам образовывали последовательность — только в таком случае, метрики *gcb* и *g-chunks* смогут учитывать их значения наиболее эффективно.

4. Метод обучения модели отрицательного отбора, использующий перестановку битов

Использование бинарных детекторов часто приводит к такому явлению, как “дыры”. “Дырой” (hole) называют такую бинарную строку, описывающую экземпляр чужого класса, для которой невозможно сгенерировать корректный детектор (т.е. любой сгенерированный детектор, который соответствует этой строке, будет также соответствовать какой-то строке, описывающей “свой” экземпляр). Иными словами, чужая строка $a \in N$ образовывает “дыру”, тогда и только тогда, когда $\forall x \in U, match(x, a) = 1: \exists s \in S, match(s, a) = 1$. “Дыры” образовываются при использовании любой метри-

ки с фиксированной вероятностью соответствия [6]. В [9] предлагается решение этой проблемы. Каждый детектор должен иметь несколько способов представления в бинарном виде. Например, для строк $s_1 = 01101011$, $s_2 = 00010011$ зададим правило перестановки бит $L = 1-6-2-5-8-3-7-4$. Применяв это правило к строкам, получим: $L(s_1) = 00111110$, $L(s_2) = 00001011$. Используя гсб-метрику с параметром $r = 3$, можно увидеть, что $match(s_1, s_2) = 1$, т.к. последние три бита этих строк совпадают. Однако $match(L(s_1), L(s_2)) = 0$.

Фактически использование маски перестановок позволяет изменять форму детектора в пространстве признаков, в то время как форма “своего” пространства остается неизменной. Каждое фиксированное представление формирует свои “дыры”, а использование нескольких представлений снизит вероятность того, что одна и та же “чужая” строка приведет к формированию “дыр” одновременно во всех представлениях детектора.

Хотя такой подход и позволяет повысить точность работы модели отрицательного отбора за счет устранения “дыр”, а также снизить количество детекторов в наборе, для каждого детектора производится несколько вычислений функции сопоставления, таким образом вычислительная сложность метода эквивалентна отрицательному отбору с цензурированием при большем количестве детекторов. Также усложняется процесс генерации набора детекторов, поскольку требуется случайным образом сгенерировать такую пару битовых строк (сам детектор и его представление после применения правила перестановки), чтобы ни одна из них не совпала ни с каким “своим” экземпляром.

Существенным недостатком метода является то, что он применим только для метрики гсб.

5. Метод обучения с маскированием

С целью устранения недостатков метода, описанного выше, авторами предлагается использовать детекторы с маскированием. Для этого необходимо расширить алфавит, на основе которого формируются детекторы, дополнив его еще одним символом – символом маски $Z : \Omega = \{0, 1, Z\}$.

Значение Z играет особую роль – оно соответствует любому значению $\{0, 1\}$ бита в битовой строке. Учитывая этот факт, существующие метрики должны быть модифицированы.

Так, метрика Хемминга для порогового значения r и битовых строк s и d длиной l может быть представлена в виде:

$$match(d, s) = \begin{cases} 1, \sum_{i=1}^l \{1 | d_i \neq Z \wedge d_i \neq s_i\} \geq r; \\ 0, \text{ в противном случае.} \end{cases}$$

Иными словами, детектор d и битовая строка s совпадают, если у них совпадают r незамаскированных битов.

Метрика гсб для маскированных детекторов вычисляется следующим образом:

$$match(d, s) = \begin{cases} 1, \exists i, i = 1 \dots l - r : \sum_{j=i}^l \{1 | d_j = Z \wedge d_j = s_j\} \geq r; \\ 0, \text{ в противном случае.} \end{cases}$$

В тех случаях, когда в результате обучения диагностической модели требуется получить набор правил, по которым будет проводиться классификация, рекомендуется применять модификацию метрики Хемминга:

$$match(d, s) = \begin{cases} 1, \sum_{i=1}^l \{1 | d_i = Z \vee d_i = s_i\} = l; \\ 0, \text{ в противном случае.} \end{cases}$$

Такая метрика предполагает, что детектор соответствует строке, если все незамаскированные биты детектора соответствуют битам строки.

Жизненный цикл детектора включает в себя следующие стадии:

- формирование полностью замаскированного детектора, т.е. такого детектора, у которого значение всех битов равно Z ;
- замена замаскированных значений битов детектора;
- добавление детектора в набор.

Ниже представлен метод обучения модели отрицательного отбора, использующей маскирование детекторов.

1. Сформировать замаскированный детектор $d = \{Z\}^n$.

2. Если $\exists s \in S : match(d, s) = 1$, тогда перейти к этапу 3, в противном случае – к этапу 5.

3. Выбрать произвольным образом бит d_i , $i = 1, \dots, l$, $d_i = Z$. Если такого бита не существует, тогда перейти к этапу 1, в противном случае – перейти к этапу 4.

4. Установить значение i -го бита детектора: $d_i = \neg s_i$. Перейти к этапу 2.

5. Добавить детектор d в набор детекторов: $D = D \cup \{d\}$. Если достигнуто достаточное для необходимого уровня точности количество детекторов, тогда перейти к этапу 6, в противном случае – перейти к этапу 1.

6. Останов.

В результате работы метода будет получен набор детекторов, у которых незамаскированные биты определяют правила, по которым можно проводить дальнейшую классификацию. Например, детектор $\{001Z1\}$ можно рассматривать следующим образом: если у экземпляра отсутствует 1-ый и 2-ой признаки, однако присутствуют 3-ий и 5-ый, то экземпляр считается чужим.

Так как метод не проверяет наличие подобных детекторов в наборе, то в процессе работы метода могут быть сгенерированы идентичные детекторы. По той же причине могут быть получены такие

пары детекторов, у которых все незамаскированные биты совпадают, например $\{010Z\}$ и $\{01ZZ\}$.

Чтобы ускорить работу обученной модели предлагается по окончании обучения удалить из набора детекторов дубликаты. Также для детекторов, отличающихся только замаскированными битами, целесообразно оставлять только те, у которых большее количество замаскированных битов. Таким образом можно сократить набор детекторов без потери точности классификации.

6. Модифицированный метод обучения с маскированием

Предложенный выше метод позволяет сократить количество детекторов в наборе, тем самым повысив скорость процесса классификации с использованием модели отрицательного отбора. Однако формирование такого набора детекторов остается ресурсоемкой задачей, поскольку каждый новый детектор после изменения некоторого бита требуется сопоставить с каждым “своим” экземпляром.

Этого можно избежать, если сохранять промежуточные детекторы, не вошедшие в состав модели, и использовать их в качестве кандидатов при генерации нового поколения детекторов.

Модифицированный метод обучения модели отрицательного отбора с маскированием состоит из следующих этапов.

1. Сформировать замаскированный детектор $c_0 = \{Z\}^n$.
2. Создать набор кандидатов в детекторы $C = \{c_0\}$.
3. Выбрать произвольным образом кандидат в детекторы из набора $c \in C$.
4. Если $\exists s \in S : match(c, s) = 1$, тогда перейти к этапу 5, в противном случае – перейти к этапу 8.
5. Выбрать произвольным образом бит c_i , $i = 1, \dots, l$, $c_i = Z$. Если такого бита не существует, тогда перейти к этапу 1, в противном случае – перейти к этапу 6.
6. Установить значение i -го бита детектора: $c_i = \neg s_i$.
7. Добавить модифицированный кандидат в набор $C = C \cup \{c\}$. Перейти к этапу 3.
8. Добавить кандидат c в набор детекторов: $D = D \cup \{c\}$. Если достигнуто достаточное для необходимого уровня точности количество детекторов – перейти к этапу 9, в противном случае – перейти к этапу 1.
9. Останов.

В первоначальной реализации метода, кандидат в детекторы состоял исключительно из замаскированных битов, а следовательно, всегда совпадал с экземплярами из набора S . Данная модификация снижает вероятность совпадения, поскольку боль-

шая часть кандидатов в детекторы уже содержит незамаскированные биты. Модифицированный метод предполагает наличие набора кандидатов в детекторы, мощность которого увеличивается по мере обучения модели.

В качестве дальнейшего развития метода можно выбирать кандидата в детекторы из набора не случайным образом, а основываясь на значении некоторой фитнес-функции, которая своим значением определяет насколько данный кандидат пригоден для формирования детекторов. Пригодность может определяться количеством замаскированных битов, а также количеством детекторов, сформированных на основе данного кандидата.

7. Эксперименты и результаты

Предложенный метод синтеза диагностической модели тестировался как на синтетических выборках, так и на практических задачах диагностирования [10] с использованием программной реализации метода на языке Python.

Синтетические тесты включали в себя наборы детекторов длиной $l = 4, \dots, 12$. Обучение проводилось на экземплярах “своего” класса, затем проводилась классификация и вычислялось значение ошибки первого рода, определяемое числом неверно распознанных чужих экземпляров выборки.

На рис. 1 представлены графики, отображающие зависимости количества детекторов в наборе D от длины битовых строк l для метода с цензурованием и для предложенной модификации метода с маскированием соответственно.

Рис. 2 отображает количество проверок соответствия детекторов и экземпляров для каждого из методов.

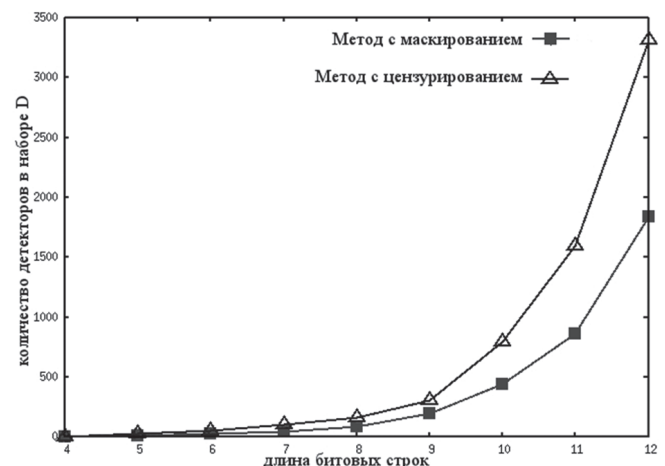


Рис. 1. График зависимости размера набора детекторов от размерности задачи

Как видно из рис. 1 и рис. 2, предложенный метод требует значительно меньше вычислений при построении модели отрицательного отбора, хотя при этом точность работы метода выше и составляет 95–100%.

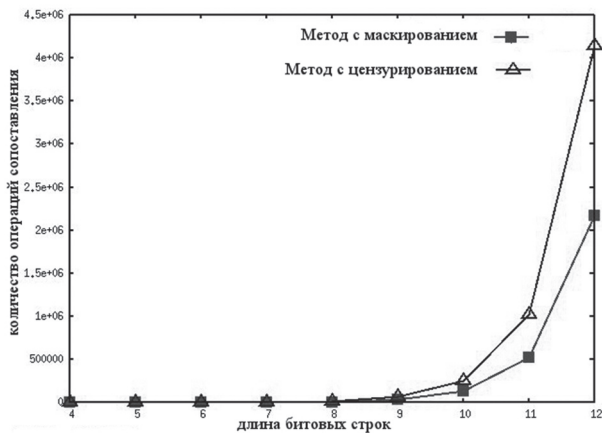


Рис. 2. График зависимости количества проверок правила сопоставления от размерности задачи

Выводы

С целью решения актуальной задачи автоматизации процесса диагностирования объектов, характеризуемых набором бинарных признаков, разработано математическое обеспечение, позволяющее проводить определение класса состояний объектов диагностирования на основе иммунокомпьютинга.

Научная новизна работы заключается в том, что впервые предложена модель отрицательного отбора, использующая маскированные детекторы, а также метод ее обучения, основная идея которого состоит в том, чтобы замаскировать максимальное количество битов, отвечающих за признаки с низкой информативностью. Это позволяет одновременно с построением диагностической модели осуществлять отбор групп информативных признаков, за счет чего повышается скорость работы метода обучения, а также упрощаются генерируемые продукционные правила посредством сокращения числа условий в антецедентах, что, в свою очередь, упрощает диагностическую модель, повышает скорость ее работы и интерпретируемость.

Практическая ценность работы заключается в том, что разработано программное обеспечение для проведения диагностирования объектов с помощью предложенной модели отрицательного отбора.

Тестирование предложенной модели отрицательного отбора показало высокую точность классификации, что позволяет рекомендовать ее использование для решения практических задач.

Дальнейшие исследования могут быть сосредоточены на развитии метода обучения предложенной модели, в частности, на разработке критериев отбора кандидатов в детекторы.

Список литературы: 1. Миркес, Е.М. Логически прозрачные нейронные сети и производство явных знаний из данных [Текст] / Е.М. Миркес, А.Н. Горбань, В.Л. Дунин-Барковский, А.Н. Кирдин и др. // Нейроинформатика.

Новосибирск: Наука. Сибирское предприятие РАН, 1998. – С. 296. 2. J. Ross Quinlan C4.5: Programs for Machine learning // Massachusetts, USA: Morgan Kaufmann Publishers, 1993. P. 324. 3. Gonzalez F., Dasgupta D., Gomez J. The effect of binary matching rules in negative selection // Proceedings of the Genetic and Evolutionary Computation Conference GECCO–2003 (9–11 July, 2003). Berlin-Heidelberg: Springer-Verlag, 2003. P. 195–206. 4. Ji Z., Dasgupta D. Revisiting negative selection algorithms // Evolutionary Computation. 2007. №15. P. 223–251. 5. Forrest S., Perelson A.S., Cherkuri R., Allen L. Self-Nonself Discrimination in a Computer // Proceedings of the 1994th IEEE Symposium on Research in Security and Privacy (1994). CA: IEEE Computer Society Press, 1994. P. 202–212. 6. D’haeseleer P., Forrest S., Helman P. An immunological approach to change detection: algorithms, analysis, and implications // Proceedings of the 1996th IEEE Symposium on Computer Security and Privacy (6–8 May, 1996). CA: IEEE Computer Society Press, 1996. P. 110–119. 7. Balthrop J., Esponda F., Forrest S., Glickman M. Coverage and generalization in an artificial immune system // Proceedings of the Genetic and Evolutionary Computation Conference GECCO–2002 (9–13 July, 2002). San Francisco, USA: Morgan Kaufmann, 2002. P. 3–10. 8. Farmer J., Packard N., Perelson A. The immune system, adaptation, and machine learning // Physica D: Nonlinear Phenomena. 1986. №2. P. 187–204. 9. Hofmeyr S., Forrest S. Architecture for an Artificial Immune System // Evolutionary Computation. 2000. №8(4). P. 443–473. 10. Герасимчук, Т.С. Использование искусственных иммунных систем для прогнозирования риска развития рекуррентных респираторных инфекций у детей раннего возраста [Текст] / Т.С. Герасимчук, С.А. Зайцев, С.А. Субботин // Матеріали міжрегіональної науково-практичної конференції “Діагностика та лікування інфекційно опосередкованих соматичних захворювань у дітей”: конференція (10–11 лютого 2011 р.). Донецьк: Норд-прес, 2011. – С. 27–29.

Поступила в редколлегию 27.06.2011

УДК 004.93

Модель негативного відбору з використанням маскованих детекторів та метод її навчання для вирішення задач діагностування / С. О. Зайцев, С. О. Субботін // Біоніка інтелекту: наук.-техн. журнал. – 2011. – № 3 (77). – С. 131–135.

Проводилось дослідження моделі негативного відбору в задачі діагностування. Запропоновано модель негативного відбору, що використовує маскування детекторів, та розроблено метод її навчання, що дозволило підвищити швидкість роботи моделі та покращити інтерпретируемість результатів. Проведено експерименти, що підтверджують ефективність запропонованої моделі.

Лл.: 2. Бібліогр.: 10 найм.

UDC 004.93

Negative selection model with masked detectods and its training method in diagnostics / S. A. Zaitsev, S. A. Subbotin // Bionics of Intellgence: Sci. Mag. – 2011. – № 3 (77). – P. 131–135.

Negative selection model application in diagnostics have been investigated. A new negative selection model based on detector masking has been proposed, which allows to increase model speed and improve results interpretation. The experiments have been carried to approve the efficiency of the suggested model.

Fig.: 2. Ref.: 10 items.