



И.В. СЕРГИЕНКО, О.А. ПРОВОТАР

УДК 681.3

**РАСПОЗНАВАНИЕ ВТОРИЧНОЙ СТРУКТУРЫ ДНК
НЕЧЕТКИМИ СИСТЕМАМИ ЛОГИЧЕСКОГО
ВЫВОДА**

Аннотация. Рассмотрены вопросы построения нечетких систем логического вывода для распознавания вторичной структуры ДНК. Приведен пример предвидения структуры центрального остатка белка MutS как выхода нечеткой системы с процедурой логического вывода Мамдани.

Ключевые слова: нечеткая система, вторичная структура ДНК, нечеткое множество.

ВВЕДЕНИЕ

Известно [1, 2], что пространственная структура ДНК определяется ее аминокислотной последовательностью. В свою очередь, пространственная структура определяет функциональность белка.

Задача распознавания структур белка различных уровней организации является достаточно сложной. Для ее решения используются различные методы и подходы, в том числе экспериментальные (основанные на физике образования химических связей), машинного обучения (используются базы данных экспериментально найденных вторичных структур как обучающих выборок), вероятностные (на основе байесовских процедур и цепей Маркова).

В настоящей статье предлагается метод распознавания вторичной структуры ДНК с помощью нечетких систем логического вывода. Задача состоит в следующем: необходимо построить нечеткую систему логического вывода, которая по произвольной аминокислотной последовательности определяла бы (в виде нечеткого множества) вторичную структуру центрального остатка (аминокислоты) входной последовательности.

Для решения этой задачи, в первую очередь, необходимо спроектировать нечеткую систему по обучающим выборкам.

ПОСТРОЕНИЕ НЕЧЕТКИХ ПРАВИЛ

Один из методов построения системы нечетких правил по числовым данным приведен в [3]. Он заключается в следующем. Пусть для простоты создается база правил с двумя входами и одним выходом. Для этого необходимы обучающие данные (выборки) в виде множества

$$(x_1(i), x_2(i); d(i)), \quad i = 1, 2, \dots, m,$$

где $x_1(i)$, $x_2(i)$ — входы модуля нечеткого управления, $d(i)$ — выход модуля нечеткого управления. Задача состоит в построении таких нечетких правил, чтобы образованная на их основе нечеткая система логического вывода по входным данным генерировала корректные выходные данные. Алгоритм решения сформулированной задачи сводится к следующей последовательности шагов.

© И.В. Сергиенко, О.А. Проватар, 2014

Шаг 1. Разделение пространства входов и выходов на области и определение соответствующих функций принадлежности. Каждый вход и выход делим на $2N + 1$ отрезок, где N — число, которое для каждого входа подбирается индивидуально. Отдельные области (отрезки) будем обозначать следующим образом:

M_N (левая N), ..., M_1 (левая 1), S (средняя), D_1 (правая 1), ..., D_N (правая N).

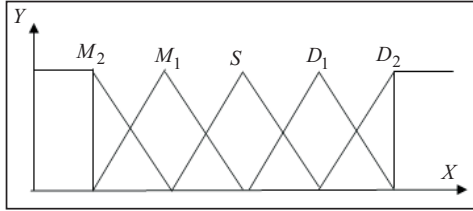


Рис. 1

Каждая функция принадлежности имеет определенную форму, например треугольную (рис. 1).

Шаг 2. Построение нечетких правил на основе обучающих выборок. Определяются меры принадлежности обучающих данных $(x_1(i), x_2(i); d(i))$ в каждой области. Для каждого набора обучающих данных нечеткое правило логического вывода формируется следующим образом:

$$\begin{aligned} (x_1(i), x_2(i); d(i)) &\Rightarrow (x_1(i) [\text{макс в соответствующей области } S_1], \\ &x_2(i) [\text{макс в соответствующей области } S_2], \\ &d(i) [\text{макс в соответствующей области } S_3]) \Rightarrow \\ &\Rightarrow R^{(i)}: \text{if } x_1 \text{ есть } S_1 \text{ and } x_2 \text{ есть } S_2 \text{ then } y \text{ есть } S_3. \end{aligned}$$

Шаг 3. Элиминация противоречий. Некоторые правила могут быть противоречивыми. Один из способов решения этой проблемы заключается в приписывании каждому правилу степени истинности с последующим выбором правила с наибольшей мерой принадлежности. Например, для правила вида

$$R: \text{if } x_1 \text{ есть } A_1 \text{ and } x_2 \text{ есть } A_2 \text{ then } y \text{ есть } B$$

мера истинности определяется как

$$SR(R) = \mu_{A_1}(x_1) \mu_{A_2}(x_2) \mu_B(y).$$

Шаг 4. Дефаззификация. Происходит при определении количественных значений выходной величины.

МОДИФИЦИРОВАННЫЙ АЛГОРИТМ ПОСТРОЕНИЯ НЕЧЕТКИХ ПРАВИЛ

Пусть создается база правил с n входами и одним выходом с использованием имеющихся обучающих данных (выборок) в виде множества пар

$$(x_1(i), x_2(i), \dots, x_n(i); d(i)), \quad i=1, 2, \dots, m,$$

где $x_j(i)$ — входы модуля нечеткого управления, $x_j(i) \in \{a_1, a_2, \dots, a_k\}$; $d(i)$ — выход модуля нечеткого управления, $d(i) \in \{b_1, b_2, \dots, b_l\}$. Следует построить на основе обучающих данных нечеткую систему логического вывода, которая по произвольным входным данным генерировала бы корректные исходные данные.

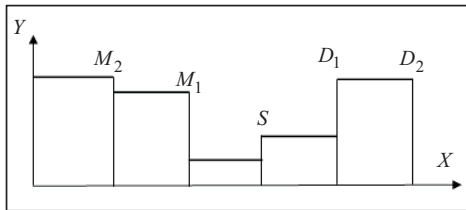


Рис. 2

Разделяем обучающие данные на группы по m_1, \dots, m_k строк, т.е. каждый вход и выход делим на $2N + 1$ отрезок, где N подбирается индивидуально для каждого входа.

Отдельные области (отрезки) будем обозначать следующим образом:

M_N (левая N), ..., M_1 (левая 1), S (средняя), D_1 (правая 1), ..., D_N (правая N).

Определим для каждой области функцию принадлежности. В случае $N = 2$ функция принадлежности может иметь вид, представленный на рис. 2.

Шаг 1. Разделение пространства входов и выходов на области. Разделяем обучающие данные на группы по m_1, \dots, m_k строк, т.е. каждый вход и выход делим на $2N + 1$ отрезок, где N подбирается индивидуально для каждого входа.

Шаг 2. Построение нечетких множеств на основе обучающих выборок.

Каждой группе m_i обучающих данных

$$(x_1(1), x_2(1), \dots, x_n(1); d(1)),$$

$$(x_1(2), x_2(2), \dots, x_n(2); d(2)),$$

$$(x_1(m_i), x_2(m_i), \dots, x_n(m_i); d(m_i))$$

сопоставляем нечеткие множества вида

$$A_1^{m_i} = \frac{|a_1^{(1)}|}{m_i} / a_1 + \dots + \frac{|a_k^{(1)}|}{m_i} / a_k,$$

.....

$$A_n^{m_i} = \frac{|a_1^{(n)}|}{m_i} / a_1 + \dots + \frac{|a_k^{(n)}|}{m_i} / a_k,$$

$$B^{m_i} = \frac{|b_1|}{m_i} / b_1 + \dots + \frac{|b_l|}{m_i} / b_l,$$

где $|a_1^{(j)}|$ — количество символов a_1 в j -м столбце группы обучающих данных, $|b_j|$ — количество символов b_j в последнем столбце группы обучающих данных.

Шаг 3. Построение нечетких правил на основе нечетких множеств.

Осуществляется на основании предыдущего шага по следующей схеме:

$$\left. \begin{array}{l} (x_1(1), x_2(1), \dots, x_n(1); d(1)) \\ (x_1(2), x_2(2), \dots, x_n(2); d(2)) \\ \dots \dots \dots \dots \dots \dots \\ (x_1(m_1), x_2(m_1), \dots, x_n(m_1); d(m_1)) \end{array} \right\} \Rightarrow$$

$$\Rightarrow R^{(1)}: \text{if } x_1 \text{ есть } A_1^{m_1} \text{ and } x_2 \text{ есть } A_2^{m_1} \text{ and } \dots \text{ and } x_n \text{ есть } A_n^{m_1} \text{ then } y \text{ есть } B^{m_1}.$$

Шаг 4. Элиминация противоречий.

Шаг 5. Дефаззификация. Проводится при определении количественных значений выходной величины.

Вышеприведенный алгоритм каждому набору обучающих данных ставит в соответствие нечеткое правило логического вывода.

Пример. Покажем, каким образом использовать предложенный алгоритм построения нечетких правил для распознавания вторичной структуры ДНК.

Известно [1, 2], что вторичная структура фрагментов полипептидной последовательности определяется в основном взаимодействиями соседних аминокислот в пределах этих фрагментов. Точнее говоря, тип вторичной структуры конкретного остатка определяется исходя из окружения соседних аминокислот.

Для построения нечетких правил логического вывода используются обучающие выборки, состоящие из 15 остатков белка MutS [2], которые имеют следующую последовательность:

<i>K</i>	<i>V</i>	<i>S</i>	<i>E</i>	<i>G</i>	<i>G</i>	<i>L</i>	I	<i>R</i>	<i>E</i>	<i>G</i>	<i>Y</i>	<i>D</i>	<i>P</i>	<i>D</i>
<i>e</i>	-	-	-	<i>h</i>	<i>h</i>	<i>h</i>	h	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>
<i>V</i>	<i>S</i>	<i>E</i>	<i>G</i>	<i>G</i>	<i>L</i>	<i>I</i>	R	<i>E</i>	<i>G</i>	<i>Y</i>	<i>D</i>	<i>P</i>	<i>D</i>	<i>L</i>
-	-	-	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	h	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>
<i>S</i>	<i>E</i>	<i>G</i>	<i>G</i>	<i>L</i>	<i>I</i>	<i>R</i>	E	<i>G</i>	<i>Y</i>	<i>D</i>	<i>P</i>	<i>D</i>	<i>L</i>	<i>D</i>
-	-	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	h	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>
<i>E</i>	<i>G</i>	<i>G</i>	<i>L</i>	<i>I</i>	<i>R</i>	<i>E</i>	G	<i>Y</i>	<i>D</i>	<i>P</i>	<i>D</i>	<i>L</i>	<i>D</i>	<i>A</i>
-	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	h	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>
<i>G</i>	<i>G</i>	<i>L</i>	<i>I</i>	<i>R</i>	<i>E</i>	<i>G</i>	Y	<i>D</i>	<i>P</i>	<i>D</i>	<i>L</i>	<i>D</i>	<i>A</i>	<i>L</i>
<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	h	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>	<i>h</i>

Распознавание относится к центральному остатку, вторичная структура h используется для спирали, e — для цилиндра, символ "-" применяется для обозначения иной структуры.

Согласно алгоритму разделяем обучающие данные, например на три группы:

$K \ V \ S \ E \ G \ G \ L \ \mathbf{h} \ R \ E \ G \ Y \ D \ P \ D$
 $V \ S \ E \ G \ G \ L \ I \ \mathbf{h} \ E \ G \ Y \ D \ P \ D \ L;$
 $S \ E \ G \ G \ L \ I \ R \ \mathbf{h} \ G \ Y \ D \ P \ D \ L \ D$
 $E \ G \ G \ L \ I \ R \ E \ \mathbf{h} \ Y \ D \ P \ D \ L \ D \ A;$
 $G \ G \ L \ I \ R \ E \ G \ \mathbf{h} \ D \ P \ D \ L \ D \ A \ L.$

Сопоставим каждой группе соответствующие нечеткие множества:

1-я группа:

$$\begin{aligned}
 A_1^{(m_1)} &= \frac{1}{2}/K + \frac{1}{2}/V, & A_6^{(m_1)} &= \frac{1}{2}/G + \frac{1}{2}/L, & A_{11}^{(m_1)} &= \frac{1}{2}/Y + \frac{1}{2}/D, \\
 A_2^{(m_1)} &= \frac{1}{2}/V + \frac{1}{2}/S, & A_7^{(m_1)} &= \frac{1}{2}/L + \frac{1}{2}/I, & A_{12}^{(m_1)} &= \frac{1}{2}/D + \frac{1}{2}/P, \\
 A_3^{(m_1)} &= \frac{1}{2}/S + \frac{1}{2}/E, & A_8^{(m_1)} &= \frac{1}{2}/R + \frac{1}{2}/E, & A_{13}^{(m_1)} &= \frac{1}{2}/P + \frac{1}{2}/D, \\
 A_4^{(m_1)} &= \frac{1}{2}/E + \frac{1}{2}/G, & A_9^{(m_1)} &= \frac{1}{2}/E + \frac{1}{2}/G, & A_{14}^{(m_1)} &= \frac{1}{2}/D + \frac{1}{2}/L. \\
 A_5^{(m_1)} &= 1/G, & A_{10}^{(m_1)} &= \frac{1}{2}/G + \frac{1}{2}/Y,
 \end{aligned}$$

2-я группа

$$\begin{aligned}
 A_1^{(m_2)} &= \frac{1}{2}/S + \frac{1}{2}/E, & A_6^{(m_2)} &= \frac{1}{2}/I + \frac{1}{2}/R, & A_{11}^{(m_2)} &= \frac{1}{2}/P + \frac{1}{2}/D, \\
 A_2^{(m_2)} &= \frac{1}{2}/E + \frac{1}{2}/G, & A_7^{(m_2)} &= \frac{1}{2}/R + \frac{1}{2}/E, & A_{12}^{(m_2)} &= \frac{1}{2}/D + \frac{1}{2}/L, \\
 A_3^{(m_2)} &= 1/G, & A_8^{(m_2)} &= \frac{1}{2}/G + \frac{1}{2}/Y, & A_{13}^{(m_2)} &= \frac{1}{2}/L + \frac{1}{2}/D, \\
 A_4^{(m_2)} &= \frac{1}{2}/G + \frac{1}{2}/L, & A_9^{(m_2)} &= \frac{1}{2}/Y + \frac{1}{2}/D, & A_{14}^{(m_2)} &= \frac{1}{2}/D + \frac{1}{2}/A. \\
 A_5^{(m_2)} &= \frac{1}{2}/L + \frac{1}{2}/I, & A_{10}^{(m_2)} &= \frac{1}{2}/D + \frac{1}{2}/P,
 \end{aligned}$$

3-я группа

$$\begin{aligned}
 A_1^{(m_3)} &= 1/G, & A_6^{(m_3)} &= 1/E, & A_{11}^{(m_3)} &= 1/L, \\
 A_2^{(m_3)} &= 1/G, & A_7^{(m_3)} &= 1/G, & A_{12}^{(m_3)} &= 1/D, \\
 A_3^{(m_3)} &= 1/L, & A_8^{(m_3)} &= 1/D, & A_{13}^{(m_3)} &= 1/A, \\
 A_4^{(m_3)} &= 1/I, & A_9^{(m_3)} &= 1/P, & A_{14}^{(m_3)} &= 1/L. \\
 A_5^{(m_3)} &= 1/R, & A_{10}^{(m_3)} &= 1/D,
 \end{aligned}$$

Выходные нечеткие множества:

$$B^{(m_1)} = 1/h, \quad B^{(m_2)} = 1/h, \quad B^{(m_3)} = 1/h.$$

Тогда система нечетких правил логического вывода будет иметь следующий вид:

- $R^{(1)}$: if x_1 есть $A_1^{(m_1)}$ and x_2 есть $A_2^{(m_1)}$ and ... and x_{14} есть $A_{14}^{(m_1)}$ then y есть $B^{(m_1)}$,
 $R^{(2)}$: if x_1 есть $A_1^{(m_2)}$ and x_2 есть $A_2^{(m_2)}$ and ... and x_{14} есть $A_{14}^{(m_2)}$ then y есть $B^{(m_2)}$,
 $R^{(3)}$: if x_1 есть $A_1^{(m_3)}$ and x_2 есть $A_2^{(m_3)}$ and ... and x_{14} есть $A_{14}^{(m_3)}$ then y есть $B^{(m_3)}$.

Существуют два основных способа определения выхода таких систем [4]. В них используется так называемое понятие агрегации правил, т.е. учета суммар-

ного эффекта от использования всех правил. Оператор агрегации **Agg** действует как s -норма [3], но допускается использование произвольной t -нормы.

Более общей является процедура, которая использует так называемые уровни истинности нечетких правил, т.е. если правило имеет вид

if x_1 есть $A_{i1} \wedge x_2$ есть $A_{i2} \wedge \dots \wedge x_n$ есть A_{in} **then** y есть B_i

(где $x_j, j=1, \dots, n$ — входные лингвистические переменные, y — выходная лингвистическая переменная; A_{ij} и B_i — нечеткие множества; логическая связка \wedge интерпретируется как t -норма нечетких множеств), то в случае двух входов x_1 и x_2 процедура выполнения алгоритма будет заключаться в выполнении следующих шагов.

1. Для каждого правила $R^{(i)}, i=1, 2, \dots, m$, вычисляем его уровень истинности

$$\alpha_i = \min [\max_{X_1} (A'_1(x_1) \wedge A_{i1}(x_1)), \max_{X_2} (A'_2(x_2) \wedge A_{i2}(x_2))].$$

2. Для каждого правила находим индивидуальные выходы

$$B'_i(y) = \min (\alpha_i, B_i(y)).$$

3. Вычисляем агрегатный выход

$$B'(y) = \max (B'_1, B'_2, \dots, B'_m).$$

Эта процедура называется \max - \min -процедурой или процедурой логического вывода Мамдани (импликация интерпретируется как операция минимум, агрегация выходов правил — как операция максимум).

Пользуясь последним алгоритмом, определим выход полученной системы нечетких правил. Пусть на вход подается следующая аминокислотная последовательность:

L K V S E G G L I R E G Y D P.

В соответствии с процедурой выполнения алгоритма будем находить уровни истинности правил логического вывода.

1. Уровень истинности первого правила:

$$\begin{aligned} \alpha_1 = \min [& \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 1 \wedge 0), \max (1 \wedge 0.5, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0.5, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0)] = 0.5. \end{aligned}$$

2. Уровень истинности второго правила:

$$\begin{aligned} \alpha_2 = \min [& \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \\ & \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0), \max (1 \wedge 0.5, 0.5 \wedge 0), \max (1 \wedge 0, 0.5 \wedge 0, 0.5 \wedge 0)] = 0.5. \end{aligned}$$

3. Уровень истинности третьего правила:

$$\begin{aligned} \alpha_3 = \min [& \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \\ & \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \\ & \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \\ & \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \\ & \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1), \max (1 \wedge 0, 0 \wedge 1)] = 1. \end{aligned}$$

Вычисляем индивидуальные выходы каждого правила:

$$B'_1(-) = \min(0.5, 0) = 0, \quad B'_1(e) = \min(0.5, 0) = 0, \quad B'_1(h) = \min(0.5, 1) = 0.5;$$

$$B'_2(-) = \min(0.5, 0) = 0, \quad B'_2(e) = \min(0.5, 0) = 0, \quad B'_2(h) = \min(0.5, 1) = 0.5;$$

$$B'_3(-) = \min(1, 0) = 0, \quad B'_3(e) = \min(1, 0) = 0, \quad B'_3(h) = \min(1, 1) = 1.$$

В результате агрегации индивидуальных выходов имеем $B' = 1/h$. Таким образом, вторичная структура остатка L есть h .

ЗАКЛЮЧЕНИЕ

С помощью рассмотренного подхода можно распознавать (предвидеть) вторичную структуру произвольных белков. Чем большее количество обучающих выборок, тем точнее будет распознавание. Учитывая, что наиболее эффективными методами распознавания вторичной структуры белка являются нейросетевые, а также их эквивалентность методам, основанным на системах нечеткого логического вывода, можно утверждать, что с помощью рассматриваемого метода можно достичь достаточно высокой точности распознавания. Кроме того, предложенный метод позволяет учитывать погрешности при определении аминокислотных остатков в обучающих выборках.

СПИСОК ЛИТЕРАТУРЫ

1. Гупал А.М., Сергиенко И.В. Оптимальные процедуры распознавания. — Киев: Наук. думка, 2008. — 382 с.
2. Леск А. Введение в биоинформатику. — М.: Лаборатория знаний, 2009. — 320 с.
3. Рутковская Д., Пилиньский М., Рутковский Л. Нейронные сети, генетические алгоритмы и нечеткие системы. — М.: Телеком, 2006. — 382 с.
4. Провотар А.И., Лапко А.В., Провотар А.А. Нечеткие системы логического вывода и их применение // Кибернетика и системный анализ. — 2013. — № 4. — С. 37–46.

Поступила 17.10.2013