

1 Построение TNFA

Для построения TNFA по дереву можно использовать несколько изменённый алгоритм Томпсона. Узлы дерева YARD можно разделить на 4 типа: лист, последовательность, альтернатива, замыкание. Эти типы соответствуют основным конструкциям расширенных регулярных выражений.

Атрибут ребра (дуги, перехода в NFA) - функция. В общем случае ей доступна вся информация о состоянии автомата и она может производить любые действия. В нашем случае она будет записывать информацию о начале и конце распознавания конструкции регулярного выражения в список.

Каждая метка соответствует началу или концу распознавания той или иной конструкции регулярного выражения и содержит:

- уникальный идентификатор, одинаковый для метки начала и конца одной и той же конструкции
- информацию о конструкции, которой она соответствует
- позицию маркера в строке, в момент генерации метки

Определим следующий набор меток:

```
type tag =  
  //Последовательность  
  | TSeqStart of int*int //(уникальный идентификатор, позиция маркера)  
  | TSeqEnd of int*int  
  //Первая ветка альтернативы  
  | TAlt1Start of int*int  
  | TAlt1End of int*int  
  //Вторая ветка альтернативы  
  | TAlt2Start of int*int  
  | TAlt2End of int*int  
  //Замыкание  
  | TClsStart of int*int  
  | TClsEnd of int*int
```

Определим функции, которые будем использовать в качестве меток для рёбер. Общим параметром для всех функций будет текущая позиция указателя в строке *pos*. Так же каждой функции доступна некоторая общая структура данных, в которую ведётся запись данных, *globalTagsList*. Набор функций выглядит следующим образом:

- $\omega - \text{fun } _ \rightarrow ()$ Будет означать отсутствие метки.
- $\text{SeqStart}_k(\text{pos}) - \text{fun pos} \rightarrow \text{TSeqStart}(k, \text{pos}) :: \text{globalTagsList}$ Подразумевается, что уникальный идентификатор вычисляется заранее, поэтому для разных конструкций строятся разные функции. Чтобы различать их, присвоим им индекс, соответствующий идентификатору.
- $\text{SeqEnd}_k(\text{pos}) - \text{fun pos} \rightarrow \text{TSeqEnd}(k, \text{pos}) :: \text{globalTagsList}$

- $Alt1Start_k(pos) - \text{fun pos} \rightarrow TAlt1Start(k, pos) :: \text{globalTagsList}$
- $Alt1End_k(pos) - \text{fun pos} \rightarrow TAlt1End(k, pos) :: \text{globalTagsList}$
- $Alt2Start_k(pos) - \text{fun pos} \rightarrow TAlt2Start(k, pos) :: \text{globalTagsList}$
- $Alt2End_k(pos) - \text{fun pos} \rightarrow TAlt2End(k, pos) :: \text{globalTagsList}$
- $ClsStart_k(pos) - \text{fun pos} \rightarrow TClsStart(k, pos) :: \text{globalTagsList}$
- $ClsEnd_k(pos) - \text{fun pos} \rightarrow TClsEnd(k, pos) :: \text{globalTagsList}$

Тогда для расстановки меток можно будет дополнить алгоритм Томпсона. Результат будет таким:

- Лист : $Leaf(a)$

$$i \xrightarrow{a/\omega} i+1$$

- Последовательность : $Seq(lst)$

$$i \xrightarrow{\epsilon/SeqStart_k(pos)} \text{concat (map doTNFA lst)} \xrightarrow{\epsilon/SeqEnd_k(pos)} i+1$$

- Альтернатива : $Alt(a,b)$

$$i \begin{array}{l} \nearrow \epsilon/Alt1Start_k(pos) \\ \searrow \epsilon/Alt2Start_k(pos) \end{array} \begin{array}{l} (\text{doTNFA } a) \\ (\text{doTNFA } b) \end{array} \begin{array}{l} \nwarrow \epsilon/Alt1End_k(pos) \\ \nearrow \epsilon/Alt2End_k(pos) \end{array} i+1$$

- Замыкание : $Cls(a)$

$$i \xrightarrow{\epsilon/ClsStart_k(pos)} i+1 \begin{array}{c} \nearrow \epsilon/\omega \\ \xrightarrow{\epsilon/\omega} \\ \nwarrow \epsilon/\omega \end{array} \begin{array}{c} (\text{doTNFA } a) \\ \nwarrow \epsilon/\omega \\ \nearrow \epsilon/\omega \end{array} i+2 \xrightarrow{\epsilon/ClsEnd_k(pos)} i+3$$

ϵ/ω

После работы построенного таким образом автомата в *globalTagsList* будет лежать структура, описывающая процесс распознавания строк. Она будет являться правильной скобочной структурой (если считать, что метки начала и конца конструкции образуют скобочную пару).

2 Построение дерева вывода

Для построения дерева вывода строки, поданной на вход автомату, можно воспользоваться деией, описанной выше. Для этого потребуется только изменить функции меток и глобальную структуру данных.

Необходимо преобразовать *globalTagsList* в стек. Теперь в нём будут лежать не только метки, но и деревья вывода. Функции *Start* по-прежнему будут класть в *globalTagsList* метку о начале соответствующей конструкции, а *End* функции будут снимать все элементы до соответствующей метки, строить из них очередной узел, соответствующий распознанной конструкции и класть его обратно.

Таким образом, в конце работы автомата, если строка принята, то в *globalTagsList* будет лежать дерево вывода этой строки.

3 Преобразование TNFA в TDFA

Общая идея преобразования заключается в том, что можно расширить понятие состояния таким образом, чтобы стал возможен перенос в него функций с рёбер. Часть функций будут выполняться на входе в состояние, часть на выходе.

4 Пример



