

*СЛОЖНОСТЬ РАСПОЗНАВАНИЯ  
ПРИНАДЛЕЖНОСТИ СЛОВА  
РЕГУЛЯРНОМУ ЯЗЫКУ, ЗАДАННОМУ  
НАБОРОМ РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ*

Александров Д.Е.

22 апреля 2014 года

# Регулярные выражения

Далее рассматриваются PCRE-совместимые регулярные выражения над языком  $\Sigma$  со следующими обозначениями:

- ▶ " $R_1 R_2$ " — конкатенация выражений  $R_1$  и  $R_2$
- ▶ " $(R_1 | R_2)$ " — объединение (ИЛИ) выражений  $R_1$  и  $R_2$
- ▶ " $R^*$ " — конкатенация произвольного числа выражений  $R$
- ▶ " $R\{m,n\}$ " — конкатенация от  $m$  до  $n$  раз выражений  $R$
- ▶ " ." — объединение всех символов языка
- ▶ " $[a_1 a_2 \dots a_n]$ " — объединение символов  $a_1, a_2 \dots a_n$
- ▶ " $[\^{} a_1 a_2 \dots a_n]$ " — объединение всех символов языка кроме символов  $a_1, a_2 \dots a_n$

# Варианты решения задачи поиска

Построение конечного автомата:

- ▶ НКА (NFA)
- ▶ ДКА (DFA)
- ▶ МДКА (MDFA)
- ▶ Д2КА ( $D^2 FA$ )
- ▶ Дуальный КА (DualFA)
- ▶ Расширенный ДКА (XFA)

# Недетерминированный конечный автомат

Заданы:

- ▶  $Q$  — множество состояний
- ▶  $A \subseteq Q$  — множество конечных состояний
- ▶  $q_0 \in Q$  — начальное состояние
- ▶  $\delta: Q \times \Sigma \rightarrow 2^Q$  — функция перехода

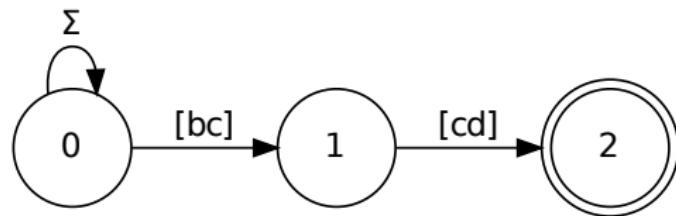
Определим функцию  $\delta^*$ :

- ▶  $\forall q \in Q, a \in \Sigma: \delta^*(q, a) = \delta(q, a)$
- ▶  $\forall q \in Q, a \in \Sigma, \alpha \in \Sigma^*: \delta^*(q, \alpha a) = \bigcup_{p \in \delta^*(q, \alpha)} \delta(p, a)$

Слово  $\alpha$  представимо НКА, если  $\delta^*(q_0, \alpha) \cap A \neq \emptyset$

# Недетерминированный конечный автомат

Пример НКА для выражения ". \* [bc] [cd]" :



$$\{0\} \xrightarrow{a} \{0\} \xrightarrow{b} \{0, 1\} \xrightarrow{c} \{0, 1, 2\} \xrightarrow{d} \{0, 2\}.$$

# *Недетерминированный конечный автомат*

Хорошо:

- ▶ малое число состояний

Плохо:

- ▶ высокая сложность вычисления

# Детерминированный конечный автомат

Заданы:

- ▶  $Q$  — множество состояний
- ▶  $A \subseteq Q$  — множество конечных состояний
- ▶  $q_0 \in Q$  — начальное состояние
- ▶  $\delta: Q \times \Sigma \rightarrow Q$  — функция перехода

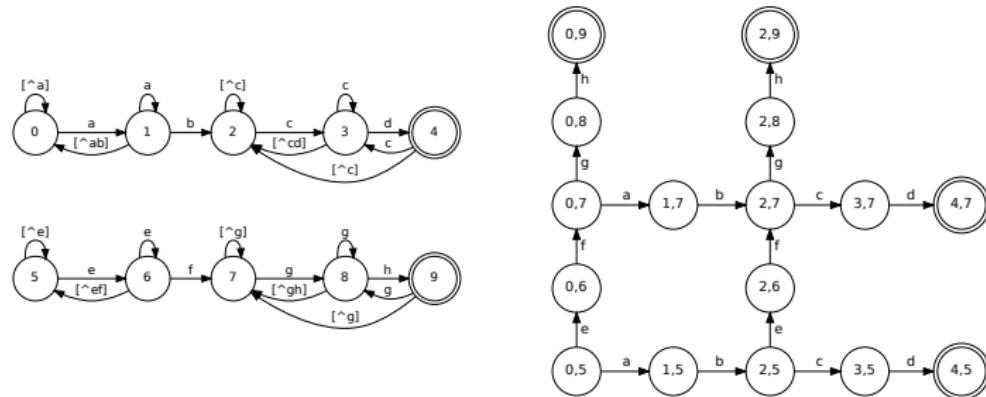
Определим функцию  $\delta^*$ :

- ▶  $\forall q \in Q, a \in \Sigma: \delta^*(q, a) = \delta(q, a)$
- ▶  $\forall q \in Q, a \in \Sigma, \alpha \in \Sigma^*: \delta^*(q, \alpha a) = \delta(\delta^*(q, \alpha), a)$

Слово  $\alpha$  представимо ДКА, если  $\delta^*(q_0, \alpha) \in A$

# Детерминированный конечный автомат

Пример ДКА для выражений “.\*ab.\*cd” и “.\*ef.\*gh”:



# Детерминированный конечный автомат

Хорошо:

- ▶ низкая сложность вычисления

Плохо:

- ▶ большое число состояний

# *МДКА*

Исходные регулярные выражения делятся на группы, по которым строятся отдельные ДКА.

Причем группировка производится так, чтобы число состояний каждого автомата не превышало заданного значения.

# *МДКА*

Хорошо:

- ▶ фиксированная сложность вычисления
- ▶ возможно небольшое число состояний

Плохо:

- ▶ возможна высокая сложность вычисления
- ▶ нетривиальное построение

*Yu F. et al. Fast and memory-efficient regular expression matching for deep packet inspection //Architecture for Networking and Communications systems, 2006. ANCS 2006. ACM/IEEE Symposium on. – IEEE, 2006. – С. 93-102.*

# Д2КА

Заданы:

- ▶  $Q$  — множество состояний
- ▶  $A \subseteq Q$  — множество конечных состояний
- ▶  $q_0 \in Q$  — начальное состояние
- ▶  $\delta: G \rightarrow Q$  — функция перехода, где  $G \subseteq Q \times \Sigma$
- ▶  $\delta^e: Q \rightarrow Q$  — функция перехода на случай “ошибки”

Алгоритм работы:

Дано: текущее состояние  $q \in Q$ , входной символ  $a \in \Sigma$

Надо: новое состояние  $q'$

$$q' := q$$

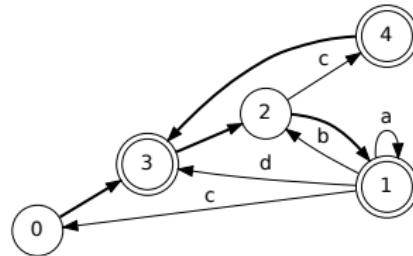
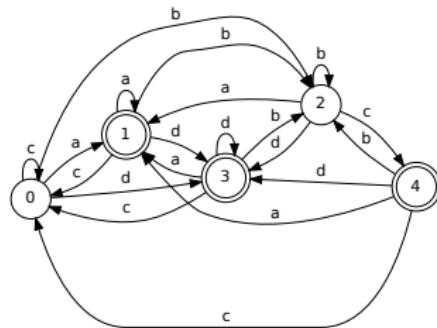
Пока  $(q', a) \notin G$

$$q' := \delta^e(q')$$

$$q' := \delta(q', a)$$

# Д2КА

Пример ДКА и Д2КА для выражений  $".*a^+"$ ,  $".*b+c"$  и  $".*c*d^+"$ :



# Д2КА

Хорошо:

- ▶ сокращение общего числа переходов

Плохо:

- ▶ возможна высокая сложность вычисления

*Kumar S. et al. Algorithms to accelerate multiple regular expressions matching for deep packet inspection //ACM SIGCOMM Computer Communication Review. – ACM, 2006. – T. 36. – №. 4. – C. 339-350.*

*Becchi M., Crowley P. An improved algorithm to accelerate regular expression evaluation //Proceedings of the 3rd ACM/IEEE Symposium on Architecture for networking and communications systems. – ACM, 2007. – C. 145-154.*

# *Дуальный автомат*

Заданы:

- ▶ ДКА с дополнительными “условными” и “расширенными” переходами
- ▶ Линейный конечный автомат — НКА, в котором из любого состояния можно попасть в не более чем 1 состояние, отличное от данного, а также не существует двух различных состояний, из которых можно перейти в третье.

## Дуальный автомат

Алгоритм работы:

Дано: ДКА в состоянии  $q \in Q$ , функция перехода  $\delta$ ,

ЛКА в состоянии  $q_I \subseteq Q_I$ , функция перехода  $\delta_I$ ,

входной символ  $a \in \Sigma$ ,

Надо: новые состояния  $q'$  и  $q'_I$

$$q' := \delta(q, a)$$

$$q'_I := \delta_I(q_I, a)$$

Если переход  $q \xrightarrow{a} q'$  — расширенный с меткой  $p \subseteq Q_I$ , то

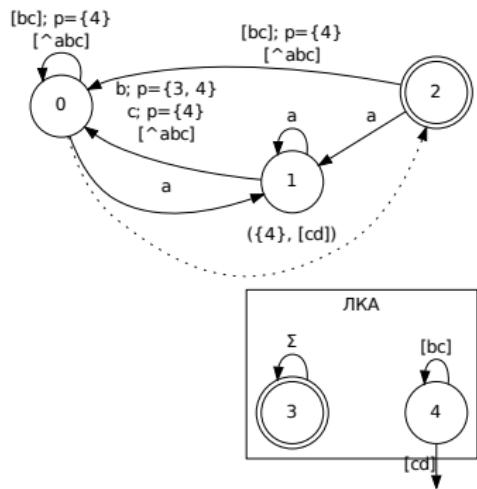
$$q'_I := q'_I \cup p$$

Если есть условный переход  $q' \xrightarrow{(q_I, a)} q''$ , то

$$q' := q''$$

# Дуальный автомат

Пример дуального автомата для выражения “.\*ab.\*” и “.\* [bc] + [cd]”:



# *Дуальный автомат*

Хорошо:

- ▶ возможно небольшое число состояний
- ▶ ЛКА реализуется через простейшие битовые операции и использует относительно немного памяти

Плохо:

- ▶ возможна высокая сложность вычисления

*Liu C., Wu J. Fast Deep Packet Inspection with a Dual Finite Automata.*  
– 2013.

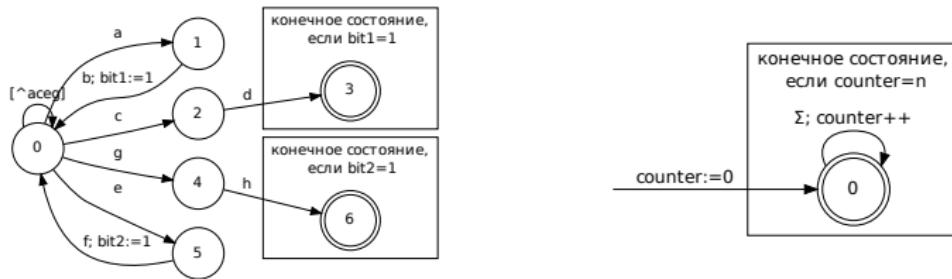
# *PДКА*

Заданы:

- ▶ набор счетчиков
  - ▶ битовых
  - ▶ числовых
- ▶ ДКА с “расширенными” переходами, изменяющими значения счетчиков

# РДКА

Пример РДКА для выражений “.\*ab.\*cd”, “.\*ef.\*gh” и “.{n}”:



# *PДКА*

Хорошо:

- ▶ возможно небольшое число состояний

Плохо:

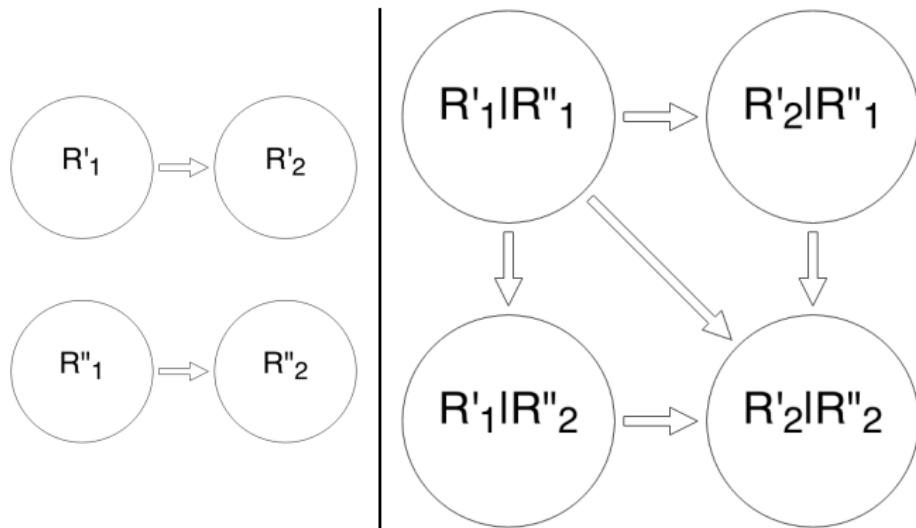
- ▶ возможна высокая сложность вычисления
- ▶ нетривиальное построение

*Smith R. et al. Deflating the big bang: fast and scalable deep packet inspection with extended finite automata //ACM SIGCOMM Computer Communication Review. – ACM, 2008. – Т. 38. – №. 4. – С. 207-218.*

# Изменение выражений

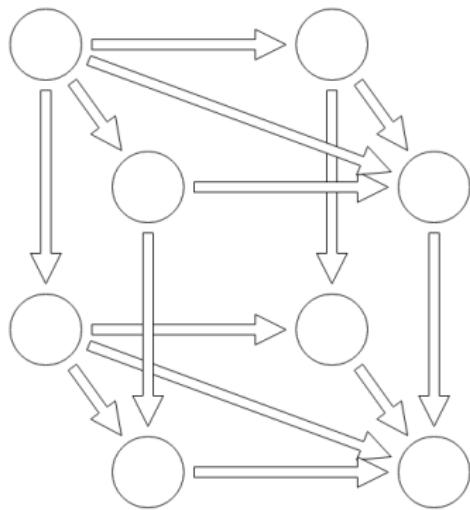
Проблема:

объединение выражений вида ". $*R'_1 * R'_2 *$ " и ". $*R''_1 * R''_2 *$ "  
ведет к быстрому росту числа состояний



# *Изменение выражений*

Схема КА для трех выражений



## Изменение выражений

Дано:

набор регулярных выражений  $R^i = . * R_1^i . * R_2^i . *$

Предложение:

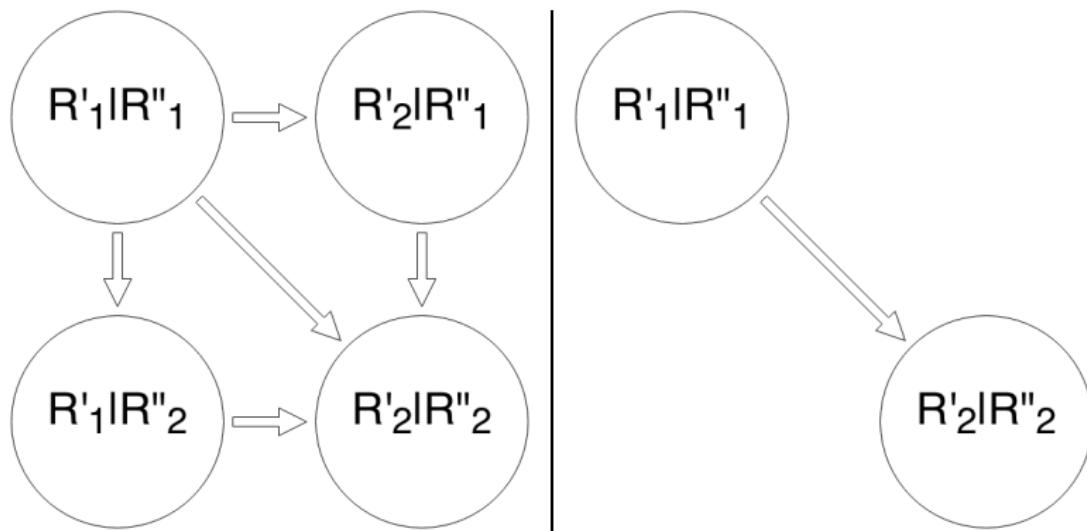
выбираем два выражения  $R^{i_1}$  и  $R^{i_2}$

заменяем их на одно:  $. * (R_1^{i_1} | R_1^{i_2}) . * (R_2^{i_1} | R_2^{i_2}) . *$

Результат:

- ▶ сокращение роста состояний КА при объединении
- ▶ ошибки при поиске (лишние слова представимы автоматом)

## *Изменение выражений*



## Изменение выражений

Вопрос:

Как можно оценивать ошибку при таком подходе?

Через  $G_i(R)$  обозначим число слов длины  $i$ , задаваемых регулярным выражением  $R$

Тогда величину ошибки можно оценить для произвольной длины  $l$ :

$$Err_l(R^1, R^2) = \frac{G_l(. * (R_1^1 | R_1^2) . * (R_2^1 | R_2^2) . *) - G_l(R^1 | R^2)}{G_l(R^1 | R^2)}$$

## Изменение выражений

Как оценить ошибку  $Err_I(R^1, R^2)$ , не вычисляя в явном виде?

Пусть далее все регулярные подвыражения  $R_i^j$  таковы, что

$$\forall \alpha \in L(R_i^j) \nexists \beta \in L(R_i^j), \beta \neq \alpha\gamma_1, \gamma_2 \in \Sigma^*: \alpha = \gamma_1\beta\gamma_2$$

Ведем функцию  $P_i(R)$ :

$$P_i(R) = \frac{G_i(R)}{|\Sigma|^i}$$

Тогда, если  $\frac{\sum_{i=0}^I P_i(R_1^1)}{\sum_{i=0}^I P_i(R_1^2)} \leq \varepsilon$  и  $\frac{\sum_{i=0}^I P_i(R_2^1)}{\sum_{i=0}^I P_i(R_2^2)} \leq \varepsilon$ , то

$$Err_I(R^1, R^2) \leq \varepsilon$$

## Изменение выражений

Кроме того, если  $\frac{\sum\limits_{i=0}^{\infty} P_i(R_1^1)}{\sum\limits_{i=0}^{\infty} P_i(R_1^2)} \leq \varepsilon$  и  $\frac{\sum\limits_{i=0}^{\infty} P_i(R_2^1)}{\sum\limits_{i=0}^{\infty} P_i(R_2^2)} \leq \varepsilon$ , то

$$Err_I(R^1, R^2) \leq \varepsilon$$