

# Теория автоматов и формальных языков

## Формальные языки и грамматики

Институт Информационных  
Технологий  
ЧелГУ, 2013

# Основные понятия

*Алфавит* – непустое конечное множество:

$$\Sigma = \{a_1, a_2, \dots, a_n\}$$

**Пример:**

$\Sigma = \{0,1\}$  - бинарный

$\Sigma = \{a,b,c,\dots,z\}$  – множество строчных букв английского алфавита

$\Sigma$  – множество всех символов (печатных символов) таблицы ASCII

*Цепочка (слово)* над алфавитом  $\Sigma$  есть конечная последовательность его элементов.

Множество всех цепочек над алфавитом  $\Sigma$  обозначим  $\Sigma^*$ .

Множество всех непустых цепочек над алфавитом  $\Sigma$  обозначим  $\Sigma^+$ .

*Длина цепочки  $x$*  есть количество её элементов и обозначается  $|x|$ .

Цепочка нулевой длины называется пустой и обозначается  $\epsilon$ .

# Основные понятия

**Пример:**

$$\Sigma = \{a, b\}$$

$$\Sigma^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$$

$$\Sigma^+ = ?$$

Цепочки ***x*** и ***y*** *равны*, если они одинаковой длины и совпадают в точности до порядка символов, из которых они состоят.

Над цепочками ***x*** и ***y*** определена операция *сцепления (произведения, конкатенации)*, заключающаяся в дописывании всех символов цепочки ***y*** после символов цепочки ***x***. Получившаяся цепочка обозначается ***xy***.

**Пример:**

***x***=кото ***y***=пёс , тогда

***xy***=котопёс ***yx***=пёското

# Основные понятия

## Понятие языка

Язык  $L$  над алфавитом  $\Sigma$  есть некоторое множество цепочек над  $\Sigma$ .

$L = \emptyset$   
 $L = \{\epsilon\}$ 
→ →
Различные языки

Формальный язык  $L$  над алфавитом  $\Sigma$  - это язык, выделенный при помощи некоторого конечного множества *формальных правил*.

Над языками определена операция *сцепления (произведения, конкатенации)*:

$$\begin{array}{l}
 L \text{ - язык} \\
 M \text{ - язык}
 \end{array}
 \quad \Rightarrow \quad
 LM = \{xy \mid x \in L, y \in M\} \text{ - язык}$$

$\{\epsilon\}L = L\{\epsilon\} = L$

# Примеры языков

множество всех слов над  $\{a, b\}$

множество  $\{a^n\}$ , где  $n$  — простое число, а  $a^n$  означает, что  $a$  повторяется  $n$  раз

множество синтаксически корректных программ в данном языке программирования

# Основные понятия

## Понятие языка – образование новых языков

- *Конкатенация (сцепление)*  $L_1L_2$  содержит все слова, удовлетворяющие форме  $vw$ , где  $v$  — это слово из  $L_1$ , а  $w$  — слово из  $L_2$ .
- *Пересечение*  $L_1 \cap L_2$  содержит все слова, содержащиеся и в  $L_1$ , и в  $L_2$ .
- *Объединение*  $L_1 \cup L_2$  содержит все слова, содержащиеся или в  $L_1$ , или в  $L_2$ .
- *Дополнение* языка  $L_1$  содержит все слова алфавита, которые не содержатся в  $L_1$ .
- *Правое отношение*  $L_1 / L_2$  содержит все слова  $v$ , для которых существует слово  $w$  в  $L_2$  такое, что  $vw$  находится в  $L_1$ .

*Пример конкатенации:*

$$L_1 = \{к, п, т\}, \quad L_2 = \{\epsilon, орт\} \quad L_1L_2 = ?$$

# Основные понятия

## Понятие языка – образование новых языков

- *Замыкание Клини*  $L_1^*$  содержит все слова, которые могут быть записаны в форме  $w_1w_2\dots w_n$ , где  $w_i$  содержится в  $L_1$  и  $n \geq 0$ . Следует помнить, что это включает и пустое слово  $\varepsilon$ , так как  $n = 0$  допустимо по условию.
- *Обращение*  $L_1^R$  содержит обращенные слова из  $L_1$ .

*Примечание:*

*обращением (или реверсом) цепочки называется цепочка, символы которой записаны в обратном порядке.*

Задача: определить замыкание Клини для  $L_1 = \{\varepsilon, 00, 1\}$ ?

# Основные понятия

$$L^* = \bigcup_{i=0}^{\infty} L^i \quad \text{- итерация}$$

$$L^+ = \bigcup_{i=1}^{\infty} L^i \quad \text{- усечённая итерация}$$

$$L^0 = \{\varepsilon\}$$

$$L^1 = L$$

$$L^{n+1} = L^n L = LL^n$$

$i$  - степень языка

**Пример:**

$$L = \{a\}$$

$$L^* = \{\varepsilon, a, aa, aaa, \dots\}$$

$$L^+ = \{a, aa, aaa, \dots\}$$

# Описание языков

## Методы описания языков

```
graph TD; A[Методы описания языков] --> B[Порождающие]; A --> C[Распознающие];
```

### Порождающие

Предполагается наличие алгоритма, последовательно порождающего все правильные предложения языка.

### Распознающие

Предполагается наличие алгоритма, определяющего принадлежность любой фразы данному языку.

# Основные понятия

## Понятие грамматики – способы задания

Простым перечислением слов, входящих в данный язык. Этот способ, в основном, применим для определения конечных языков и языков простой структуры.

Словами, порождёнными некоторой формальной грамматикой (иерархия Хомского)

Словами, порождёнными регулярным выражением

Словами, распознаваемыми некоторым конечным автоматом

Словами, порождёнными БНФ-конструкцией (Форма Бэкуса — Наура)

# Основные понятия

Понятие грамматики – способы задания

Форма Бэкуса — Наура

<Терминал>:=Задание\_терминала

<...> - ограничивают (выделяют) терминалы

| - разделяет несколько альтернативных терминалов

[...] – конструкция встречается 1 раз или отсутствует

{...} – конструкция повторяется некоторое (возможно нулевое) количество раз

Определение идентификатора:

<Буква> := A|B|C|...|Z

<Цифра> := 0|1| ... |9

<Идент> := <Буква> {<Буква>|<Цифра>}

A B A1 B2C3 DA123

<http://kitaev2005.narod.ru/tr1.htm>

# Основные понятия

Понятие грамматики – способы задания

Форма Бэкуса — Наура

Грамматика целых чисел без знака:

$\langle \text{число} \rangle := \langle \text{цифра} \rangle | \langle \text{цифра} \rangle \langle \text{число} \rangle$

$\langle \text{цифра} \rangle := 0 | 1 | 2 | \dots | 9$

Формула с плюсами и минусами без скобок.

$\langle \text{формула} \rangle := \langle \text{число} \rangle | \langle \text{формула} \rangle \langle \text{знак} \rangle \langle \text{число} \rangle$

$\langle \text{знак} \rangle := + | -$

Задачи:

1. Описать терминал  $\langle \text{формула} \rangle$  (с +, -, \*, / [и круглыми скобками])
2. Описать оператор while, а затем if (с некоторой степенью упрощения).

<http://kitaev2005.narod.ru/tr1.htm>

# Основные понятия

Порождающая грамматика – это четвёрка:

$$G = (V_T, V_N, P, S)$$

$V_T$  - конечный алфавит, определяющий множество терминальных символов

$V_N$  - конечный алфавит, определяющий множество нетерминальных символов

$P$  - конечное множество правил вывода (продукций), т.е. правил вида:

$$u \rightarrow v \\ u, v \in (V_T \cup V_N)^*$$

$S$  - начальный символ (аксиома грамматики)

$$S \in V_N$$

$$\begin{array}{|l} V = V_T \cup V_N \\ V_T \cap V_N = \emptyset \end{array}$$

# Основные понятия

*Аксиомой* называется символ в левой части первого правила вывода грамматики.

Цепочки языка, порождённого грамматикой, состоят из терминальных символов.

Будем обозначать терминальные символы - строчными, а нетерминальные символы – заглавными буквами.

В грамматике  $G$  цепочка  $x$  порождает цепочку  $y$  (обозначается  $x \Rightarrow y$ ), если:

$$\left| \begin{array}{l} x = \alpha u \beta \\ y = \alpha v \beta \\ u \rightarrow v \end{array} \right.$$

В этом случае также говорят, что  $y$  выводится из  $x$ .

Языком, порождаемым грамматикой  $G$ , называется множество терминальных цепочек, выводимых из некоторой аксиомы:

$$L(G) = \{x \mid x \in V_T^*; S \Rightarrow x\}$$

# Пример определения языка

## Задача:

$$G = (V_T, V_N, P, S)$$

$$V_T = \{a, b\}$$

$$V_N = \{S\}$$

$$P = \{S \rightarrow aSb, S \rightarrow ab\}$$

Определить язык, определяемый данной грамматикой.

## Решение:

$$S \rightarrow ab$$

$$S \rightarrow aSb \Rightarrow aabb$$

$$S \rightarrow aSb \Rightarrow aaSbb \Rightarrow aaabbb$$

...

Таким образом:

$$L(G) = \{a^n b^n \mid n \in \mathbb{N}\}$$

# Классификация Хомского

*Грамматика типа 0:*

На правила вывода не накладывается никаких ограничений.

*Грамматика типа 1 (контекстные грамматики, контекстно-зависимые):*

Все правила таких грамматик имеют вид:

$$xAy \rightarrow x\varphi y, \quad A \in V_N, \quad x, y \in V^*, \quad \varphi \in V^+$$

(Прим. Если  $\alpha \rightarrow \beta$ , то  $|\alpha| \leq |\beta|$ )

*Грамматика типа 2 (контекстно-свободные грамматики):*

Все правила таких грамматик имеют вид:

$$A \rightarrow \varphi, \quad A \in V_N, \quad \varphi \in V^+$$

*Грамматика типа 3 (автоматные грамматики, регулярные):*

Леворекурсивные содержат только правила вида:

$$A \rightarrow Aa, A \rightarrow a, \quad A \in V_N, \quad a \in (V_T \cup V_N)^*$$

Праворекурсивные содержат только правила вида:

$$A \rightarrow aA, A \rightarrow a, \quad A \in V_N, \quad a \in (V_T \cup V_N)^*$$

# Пример 1

Дана грамматика, выписать все порождаемые её цепочки и определить длину их вывода:

---

$$G = (V_T, V_N, P, S)$$

$$V_T = \{a, d, e\}$$

$$V_N = \{B, C, S\}$$

$$P = \{S \rightarrow aB, B \rightarrow Cd, B \rightarrow dC, C \rightarrow e\}$$

---

# Пример 1

Дана грамматика, выписать все порождаемые её цепочки и определить длину их вывода:

---

$$G = (V_T, V_N, P, S)$$

$$V_T = \{a, d, e\}$$

$$V_N = \{B, C, S\}$$

$$P = \{S \rightarrow aB, B \rightarrow Cd, B \rightarrow dC, C \rightarrow e\}$$

---

$$S \rightarrow aB \rightarrow aCd \rightarrow aed$$

$$S \rightarrow aB \rightarrow adC \rightarrow ade$$

Длина вывода всех (обеих, их всего 2) цепочек равна 3.

# Пример 2

Дана грамматика, определить, принадлежит ли некоторая цепочка порождаемому её языку:

---

$$G = (V_T, V_N, P, C)$$

$$V_T = \{a, b, c, d, e\}$$

$$V_N = \{A, B, C, D, E\}$$

$$P = \{A \rightarrow ed, B \rightarrow Ab, C \rightarrow Bc, C \rightarrow dD, D \rightarrow Ae, E \rightarrow bc\}$$

Цепочка для проверки: *eadabcb*

---

# Пример 2

Дана грамматика, определить, принадлежит ли некоторая цепочка порожаемому её языку:

---

$$G = (V_T, V_N, P, C)$$

$$V_T = \{a, b, c, d, e\}$$

$$V_N = \{A, B, C, D, E\}$$

$$P = \{A \rightarrow ed, B \rightarrow Ab, C \rightarrow Bc, C \rightarrow dD, D \rightarrow Ae, E \rightarrow bc\}$$

Цепочка для проверки: *eadabcb*

---

Выведем всевозможные цепочки из указанных правил вывода:

$$C \rightarrow Bc \rightarrow Abc \rightarrow edbc$$

$$C \rightarrow dD \rightarrow dAe \rightarrow dede$$

Указанную цепочку вывести никак не получается, следовательно, она не принадлежит порождённому указанной грамматикой языку.

# Пример 3

Построить контекстно-свободную грамматику, порождающую цепочки из 0 и 1 с одинаковым количеством и тех и других.

---

$$G = (V_T, V_N, P, S)$$

$$V_T = \{0, 1\}$$

$$V_N = \{S\}$$

# Пример 3

Построить контекстно-свободную грамматику, порождающую цепочки из 0 и 1 с одинаковым количеством и тех и других.

---

$$G = (V_T, V_N, P, S)$$

$$V_T = \{0, 1\}$$

$$V_N = \{S\}$$

$$P = \{S \rightarrow \varepsilon, S \rightarrow 0S1, S \rightarrow 1S0, S \rightarrow 10S, S \rightarrow 01S, S \rightarrow S10, S \rightarrow S01\}$$

# Пример 4 и 5

4. Описать грамматику языка, содержащего палиндромы. Определить тип грамматики.
  5. Описать грамматику, которая представляет выражения с операциями + и \*. (Пусть в образовании идентификаторов переменных участвуют только  $0, 1, a, b$ ). Определить тип грамматики.
-

# Пример 4 и 5

4. Описать грамматику языка, содержащего палиндромы. Определить тип грамматики.

5. Описать грамматику, которая представляет выражения с операциями + и \*. (Пусть в образовании идентификаторов переменных участвуют только 0,1,a,b). Определить тип грамматики.

Продукции для  
грамматики 4

$$P \rightarrow \varepsilon$$

$$P \rightarrow 0$$

$$P \rightarrow 1$$

$$P \rightarrow 0P0$$

$$P \rightarrow 1P1$$

Продукции для  
грамматики 5

$$E \rightarrow I$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$I \rightarrow a$$

$$I \rightarrow b$$

$$I \rightarrow Ia$$

$$I \rightarrow Ib$$

$$I \rightarrow I0$$

$$I \rightarrow I1$$

# Регулярные выражения

*Регулярные выражения задают языки.*

Обозн. Пусть  $E$  - регулярное выражение,  $L(E)$  – язык, заданный этим выражением.

Пример:  $01^*+10^*$

*Базис при построении регулярных выражений:*

1. Константы  $\epsilon$  и  $\emptyset$ . Определяют языки  $L(\epsilon)=\{\epsilon\}$  и  $L(\emptyset)=\{\emptyset\}$ .
2. Если  $a$  – произвольный символ алфавита, то  $\mathbf{a}$  – регулярное выражение, определяющее язык  $L(\mathbf{a})=\{a\}$ .
3. Переменная, обозначенная прописной буквой, например,  $L$ , представляет язык.

# Регулярные выражения

*Индукция (операторы и скобки):*

1. Если  $E$  и  $F$  — регулярные выражения, то  $E + F$  — регулярное выражение, определяющее объединение языков  $L(E)$  и  $L(F)$ , т.е.  $L(E + F) = L(E) \cup L(F)$ .
2. Если  $E$  и  $F$  — регулярные выражения, то  $EF$  — регулярное выражение, определяющее конкатенацию языков  $L(E)$  и  $L(F)$ . Таким образом,  $L(EF) = L(E)L(F)$ .
3. Если  $E$  — регулярное выражение, то  $E^*$  — регулярное выражение, определяющее итерацию языка  $L(E)$ . Таким образом,  $L(E^*) = (L(E))^*$ .
4. Если  $E$  — регулярное выражение, то  $(E)$  — регулярное выражение, определяющее тот же язык  $L(E)$ , что и выражение  $E$ . Формально,  $L((E)) = L(E)$ .

# Пример 1

Написать регулярное выражение для множества цепочек, состоящих из чередующихся нулей и единиц.

# Пример 1

Написать регулярное выражение для множества цепочек, состоящих из чередующихся нулей и единиц.

Первый вариант решения:

$$(01)^* + (10)^* + 0(10)^* + 1(01)^*$$

# Пример 1

Написать регулярное выражение для множества цепочек, состоящих из чередующихся нулей и единиц.

Первый вариант решения:

$$(01)^* + (10)^* + 0(10)^* + 1(01)^*$$

А если немного подумать, то получится и второй вариант решения:

$$(\epsilon + 1)(01)^*(\epsilon + 0)$$

# Задачи

Напишите регулярные выражения для следующих языков:

- а) (\*) множество цепочек с алфавитом  $\{a, b, c\}$ , содержащих хотя бы один символ  $a$  и хотя бы один символ  $b$ ;
- б) множество цепочек из нулей и единиц, в которых десятый от правого края символ равен 1;
- в) множество цепочек из нулей и единиц, содержащих не более одной пары последовательных единиц.

(!) Напишите регулярные выражения для следующих языков:

- а) (\*) множество всех цепочек из нулей и единиц, в которых каждая пара смежных нулей находится перед парой смежных единиц;
- б) множество цепочек, состоящих из нулей и единиц, в которых число нулей кратно пяти.