

ТЕОРЕТИЧЕСКАЯ ИНФОРМАТИКА

Алфавиты, слова, языки,
алгоритмические проблемы

Кафедра информатики и вычислительной техники

Доцент, к.т.н. Дамов Михаил Витальевич

Цели и задачи

1. Ввести подходящий формализм для работы с текстами - представлениями данных.
 - Основные понятия - алфавит, слово и язык.
2. Показать, как использовать введённые понятия, применять их для получения формальных представлений алгоритмических проблем.
 - Проблемы принадлежности (разрешимости)
 - Оптимизационные проблемы
3. Рассмотреть некоторые вопросы, связанные со сжатием текстов.
 - Сложность по Колмогорову

Формальный язык

Алфавитом называется любое непустое конечное множество.

Каждый элемент алфавита называется символом.

Алфавит языка – множество символов (букв)

Язык – множество строк

Строка (слово) – последовательность символов

Примеры:

“студент”, “123”, “house”

$\Sigma = \{ '0', '1', '2', '3', '4', '5', '6', '7', '8', '9' \}$

$\Sigma = \{ 'a', 'b', 'c', \dots, 'z' \}$ $\Sigma = \{ 'a', 'б', 'в', \dots, 'я' \}$

Примеры стандартных алфавитов

$\Sigma_{\text{bool}} = \{0, 1\}$ - логический (Булевый) алфавит

$\Sigma_{\text{lat}} = \{a, b, c, \dots, z\}$ - латинский алфавит

$\Sigma_{\text{keyboard}} = \Sigma_{\text{lat}} \cup \{\dots\}$ – алфавит символов, которые можно набрать на клавиатуре

$\Sigma_m = \{0, 1, \dots, m-1\}$, $m > 1$ – алфавит для записи чисел в m -ичной системе счисления

$\Sigma_{\text{logic}} = \{0, 1, x, (,), \cap, \cup, \neg\}$ – алфавит формул алгебры логики

Алфавит и строки

Будем использовать алфавит из двух букв

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

Строки (слова)

a, ab, abba, baba, aaabbbaabab

u = ab

v = bbbaaa

w = abba

Операции над строками

$$w = a_1 a_2 \dots a_n$$

$$v = b_1 b_2 \dots b_m$$

$$w = abba$$

$$v = bbbaaa$$

Конкатенация:

Длина строки:

$$wv = a_1 a_2 \dots a_n b_1 b_2 \dots b_m$$

$$|w| = m$$

$$wv = abbabbbaaa$$

$$|abba| = 4$$

Обращение:

Длина конкатенации строк:

$$v^R = b_m \dots b_2 b_1$$

$$|uv| = |u| + |v|$$

$$v^R = aaabbb$$

Операции над строками

Пустая строка:

Строка, не содержащая букв: λ $|\lambda| = 0$

$$\lambda w = w\lambda = w$$

$$\lambda abba = abba\lambda = abba$$

Подстрока:

Строка: $abbab$

Подстроки: $ab, abba, b, bbab, \lambda$

Операции над строками

Префиксы и суффиксы:

Строка: abba

Префиксы:

λ abba

a bba

ab ba

abb a

abba λ

Суффиксы:

u - префикс

v - суффикс

$w = uv$

Операции над строками

Итерация: $w^n = \underbrace{ww \cdots w}_n = \lambda$

$$(abba)^2 = abbaabba = ab^2a^2b^2a \quad (abba)^0 = \lambda$$

Звезда Клини:

Σ^* - множество всех возможных слов в алфавите Σ

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

$$\Sigma^* = \{\lambda, a, b, aa, ab, ba, bb, \dots\}$$

Операции над строками

Плюс Клини:

$$\Sigma^+ = \Sigma^* - \lambda \qquad \Sigma^+ = \{a, b, aa, ab, ba, bb, \dots\}$$

Язык в алфавите Σ - любое подмножество Σ^*

Операции над языками

Обычные теоретико-множественные операции:

Объединение: $\{a, ab, aaaa\} \cup \{bb, ab\} = \{a, ab, bb, aaaa\}$

Пересечение: $\{a, ab, aaaa\} \cap \{bb, ab\} = \{ab\}$

Разность: $\{a, ab, aaaa\} - \{bb, ab\} = \{a, aaaa\}$

Дополнение: $\bar{L} = \Sigma^* - L$

$\overline{\{a, ba\}} = \{\lambda, b, aa, ab, ba, bb, aaa, \dots\}$

Операции над языками

Обращение

$$L^R = \{w^R : w \in L\}$$

$$\{ab, aab, baba\}^R = \{ba, baa, abab\}$$

Конкатенация

$$L_1 L_2 = \{xy : x \in L_1, y \in L_2\}$$

$$\{a, ab, ba\}\{b, aa\} = \{ab, aaa, abb, abaa, bab, baaa\}$$

Операции над языками

Итерация

$$L^n = \underbrace{LL \dots L}_n$$

$$\{a, b\}^3 = \{a, b\} \{a, b\} \{a, b\} = \{aaa, aab, aba, abb, baa, bab, bba, bbb\}$$

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

$$\{a, b\}^0 = \{\lambda\}$$

Операции над языками

Звезда Клини (замыкание)

$$L^* = L^0 \cup L^1 \cup L^2 \cup \dots$$

$$\{a, bb\}^* = \left\{ \begin{array}{l} \lambda \\ a, bb \\ aa, abb, bba, bbbb \\ aaa, aabb, abba, bbbb, \dots \end{array} \right\}$$

Операции над языками

Плюс Клини (положительное замыкание)

$$L^+ = L^1 \cup L^2 \cup \dots = L^* - \{\lambda\}$$

$$\{a, bb\}^+ = \left\{ \begin{array}{l} a, bb \\ aa, abb, bba, bbbb \\ aaa, aabb, abba, bbbb, \dots \end{array} \right\}$$

Грамматики

Грамматики определяют языки, является ли данное предложение правильным предложением данного языка.

Пример: грамматика русского языка

<предложение> → <подлежащее> <сказуемое> <дополнение>

<подлежащее> → <существительное>

<сказуемое> → <глагол>

<дополнение> → <наречие>

<существительное> → птица | студент

<сказуемое> → летает | учится

<наречие> → высоко | хорошо

Вывод предложения

<предложение> => <подлежащее> <сказуемое> <дополнение>
=>

=> <существительное> <глагол> <наречие> =>

=> Птица летает высоко

Возможные предложения

Птица летает хорошо

Птица учится высоко

Студент летает хорошо

Обозначения

<глагол> → *летает*

<глагол> → *учится*

Переменная
или
Нетерминальный
символ

Правило
вывода

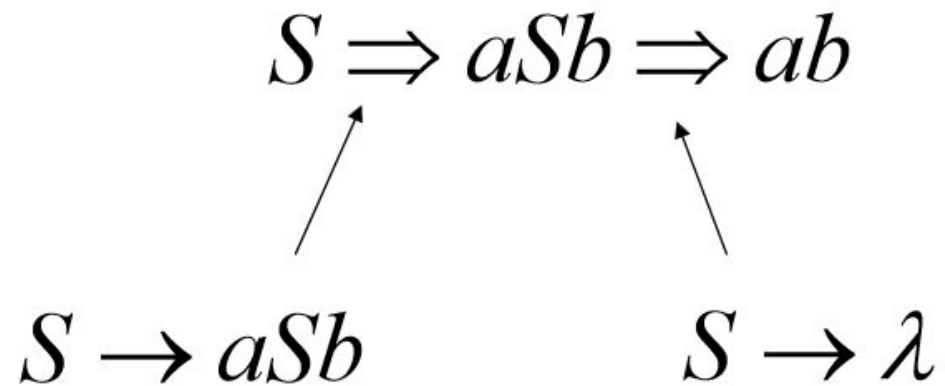
Терминальный
символ

Пример формальной грамматики

Грамматика: $S \rightarrow aSb$

$S \rightarrow \lambda$

Вывод предложения ab :



Определение формальной грамматики

$$G = \{V, T, S, P\}$$

V = Множество нетерминальных символов

T = Множество терминальных символов

S = Начальный символ

P = Множество правил вывода (продукций)

Пример:

Грамматика $S \rightarrow aSb, S \rightarrow \lambda$

$$V = \{S\}$$

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

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

Язык, порождаемый грамматикой

Определение:

Для грамматики G с начальным символом S

$$L(G) = \{w: S \Rightarrow w\}$$

-язык, порождаемый этой грамматикой

Пример:

Грамматика $G \{S \rightarrow Ab, A \rightarrow aAb, A \rightarrow \lambda\}$

$$L(G) = \{a^n b^n b: n \geq 0\}$$

поскольку $S \Rightarrow a^n b^n b$ и никакие другие слова не выводимы

Алгоритмические проблемы

Любая программа (алгоритм) A выполняет отображение:

$$A: \Sigma_1^* \rightarrow \Sigma_2^*$$

- входы представлены как слова над алфавитом Σ_1
- выходы представлены как слова над алфавитом Σ_2
- A однозначно определяет выход по каждому входу

Для некоторого алгоритма A и входа x обозначим записью $A(x)$ выход алгоритма A для этого входа. Будем говорить, что два алгоритма (программы) A и B эквивалентны, если они работают над одним и тем же алфавитом Σ и при этом $A(x) = B(x)$ для всех $x \in \Sigma^*$.

Проблема принадлежности

Обычные теоретико-множественные операции:

Объединение: $\{a, ab, aaaa\} \cup \{bb, ab\} = \{a, ab, bb, aaaa\}$

Пересечение: $\{a, ab, aaaa\} \cap \{bb, ab\} = \{ab\}$

Разность: $\{a, ab, aaaa\} - \{bb, ab\} = \{a, aaaa\}$

Дополнение: $\bar{L} = \Sigma^* - L$

$$\overline{\{a, ba\}} = \{\lambda, b, aa, ab, ba, bb, aaa, \dots\}$$
$$A(x) = \begin{cases} 1, & \text{если } x \in L \\ 0, & \text{если } x \notin L \end{cases}$$

Оптимизационная проблема

$U = \{\Sigma_I, \Sigma_O, L, M, \text{cost}, \text{goal}\}$

Σ_I – входной алфавит

Σ_O – выходной алфавит

L – язык подходящих входов

M – множество допустимых решений

Cost – функция стоимости

Goal – цель (максимизация или минимизация)

Пример: задача коммивояжера

Задание на лабораторную работу

Написать программу на языке C++, в которой необходимо выполнить следующие операции:

1. Задать алфавит языка (3 – 5 латинских букв)
2. Задать максимально возможную длину слова (5 – 7 символов)
3. Построить словарь языка.
4. Используя грамматику слайда 16, отнести случайным образом слова к трем классам.
5. Сгенерировать текст из заданного количества случайных предложений.
6. Сжать текст, используя операцию итерации.
7. Вывести в отдельные файлы: словарь, текст, сжатый текст