

1

Курс простой. По нему зачёт.

Нужно решить 4 лабораторные работы:

- Регулярные выражения в языке **Perl**. Самостоятельное изучение материала и сдача в PCMS. Если сдавать в 2022, то можно решить N-3 задач. До конца сессии -2. После -1 (все, кроме одной).
- До 22:00 в пн нужно будет записываться на сдачу лабораторных 2,3,4. Показываешь, что работает. Получаешь модификацию (их много разных на одно и то же задание). Потом показываешь уже модификацию.
- ЛР2 – ручное построение нисходящих парсеров
- ЛР3 – использование автоматических генераторов парсеров. ANTLR, Нарру, Bison (YACC)
- ЛР4 – написать автоматический генератор парсеров

Компиляция: программа → парсинг → генерация исполнимого кода.

Здесь пройдем первый основы парсинга (довольно глубоко), а генерацию кода только базовую основу.

Если вы пойдёте в магистратуру, там будет курс компиляторов. Там фокус наоборот на генерации.

В курсе нет дедлайнов. НО удачи вам с тайм-менеджментом без них :v) Курс по софт-скиллам, you're welcom

Что такое пасрер?

текстовое представление информации → внутреннее представление компьютера. Чаще всего дерево разбора в некоторой контекстно-свободной грамматике.

Регулярное выражение – способ записать регулярные языки.

Контекстно свободная грамматика – алфавит, множество нетерминалов, стартовый нетерминал и правила.

неоднозначная:

$S \rightarrow (S)$

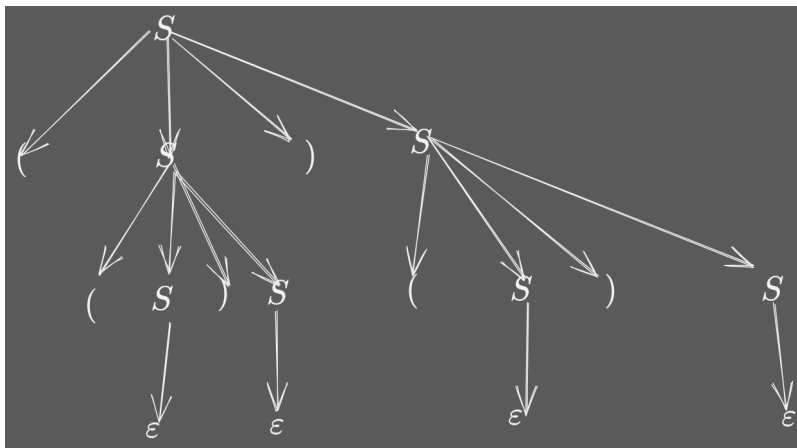
$S \rightarrow SS$

$S \rightarrow e$

однозначная:

$S \rightarrow (S)S$

$S \rightarrow e$



Токенизация: инпут в *любом* формате $\rightarrow c_1, c_2, \dots, c_n$, где c_i это токены.

Синтаксический/Лексический анализ (парсинг)

Лексический анализатор:

Есть входной алфавит парсера. Есть совсем входной алфавит A (ASCII условно).

Пример: арифметические выражения $\Sigma = \{+, -, *, (,), number\}$. Числа здесь не нужны. 2+3, 5+5, 17+231. Всё это считается одним выражением. Поэтому есть один токен для числа.

```
+ +
- -
* *
) )
( (
n (1|2|3|4|5|6|7|8|9)(0|1|2| ... |9)*|0

17+231
n
n
n+
n+n
n+ n
n+ n
```

Задача: у вас есть ЯП довольно мощный. Вы прочитали токен. Вам нужно выдать номер токена. Чтобы делать это быстро нужно Perfect < ... > Mapping

Парсеры

Алгоритм Кока-Янгера-Касами

Работает только для грамматик в нормальной форме Хомского ($A \rightarrow BC$). Можно снять ограничение, но это неприятно. Работает за куб от длины слова. Если бы мы парсили за куб ... было бы очень грустно

`dp_A[l][r]` – можно ли получить `s[l..r]` из `A`

Парсить хочется быстро.

Снизу вверх vs Сверху вниз.

Направление построения дерева. Его можно строить от корня или от листьев (спрашивая какой у них может быть родитель)

Рекурсивный парсер сверху вниз = рекурсивный спуск.

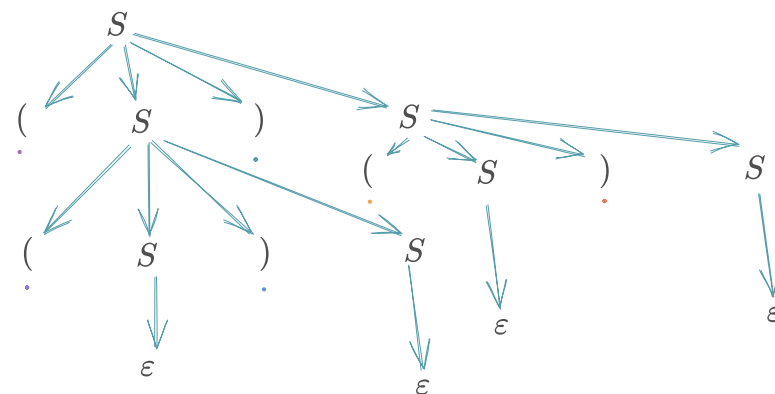
Есть языки, которые так не распарсить.

^ это было введение

Нисходящие методы разбора

$$S \rightarrow (S)S$$
$$S \rightarrow e$$

(())(())



LL(1) грамматика – знание первого нераскрытого терминала и первого терминала даёт понять следующий шаг разбора.

$$S \implies {}^*xA\alpha \implies x\xi\alpha \implies {}^*xc\alpha' \implies {}^*xcy^*$$
$$S \implies {}^*xA\beta \implies x\eta\beta \implies {}^*xc\beta' \implies {}^*xcz*$$

LL(1) грамматика – Если есть два таких вывода, то $\xi = \eta$

Патч: все символы уже подвесили, а первый нераскрытый терминал всё ещё есть. Нужно раскрывать в ε ..

Добавим символ конец ввода \$

Сейчас определение выше рассматривает бесконечное число выражений (импликаций). Можно сделать, чтобы конечное, об этом позже.

2

$\text{FIRST} : (N \cup \Sigma) \rightarrow \mathbb{P}(\Sigma \cup \varepsilon)$

$\text{FIRST}(\alpha) = \{c \mid \alpha \Rightarrow c\beta\} \cup \{\varepsilon \mid \alpha \Rightarrow \varepsilon\}$

Лемма:

- $\text{FIRST}(\varepsilon) = \{\varepsilon\}$
- $\text{FIRST}(c\xi) = \{c\} \quad c \in \Sigma$
- $\text{FIRST}(A\xi) = (\text{First}[A] \setminus \{\varepsilon\}) \cup \{\text{First}(\xi) \mid \varepsilon \in \text{First}[A]\}$

Алгоритм 1:

$\text{FIRST} = \{A \Rightarrow \emptyset \text{ for } A \in N\}$

```
while changes
  for A → alpha \in Γ
    FIRST[A] U= FIRST(alpha)
```

Note: 2126 (140) – прыскалка для доски

Лемма 2: Алгоритм 1 корректно строит First[]

Пусть $c \in \text{First}(A) \setminus \text{First}[A]$. $A \Rightarrow c\xi$ за k шагов (минимальное k)

- $k \neq 1$
- $k > 1 \quad A \Rightarrow B\eta \xrightarrow[k-1]{} c\xi$

$\text{FOLLOW} : N \rightarrow \mathbb{P}(\Sigma \cup \{\$\})$

$\text{FOLLOW}(A) = \{c \mid S \xRightarrow{*} \alpha A c \beta\} \cup \{\$ \mid S \xRightarrow{*} \alpha A\}$

Алгоритм 2:

$\text{FOLLOW} = \{A \Rightarrow \emptyset \text{ for } A \in N\}$

```
while changes:
  for A → alpha \in Γ:
    for B: alpha = xi N eta:
      FOLLOW[B] U= (FIRST[eta] \ epsilon) U
                  (FOLLOW[A] if epsilon \in FIRST(eta))
```

#png картинка дерева разбора

Пример:

$E \rightarrow T + E$ – первое слагаемое к сумме всего остального
 \wedge это не то что мы хотим

что мы хотим:

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow T * F$

$T \rightarrow F$

$F \rightarrow \text{eta}$

$F \rightarrow (E)$

$\Sigma = \{ (,), +, *, \eta \}$

$N = \{ E, T, F \}$

$FIRST[E] = FIRST[T] = FIRST[F] = \eta($

	<i>FIRST</i>	<i>FOLLOW</i>
<i>E</i>	$\eta($	$\$+$
<i>T</i>	$\eta($	$\$ + *$
<i>F</i>	$\eta($	$\$ + *$

Теорема: грамматика Γ является LL1 тогда и только тогда, когда:

- $A \rightarrow \alpha, A \rightarrow \beta \implies FIRST(\alpha) \cap FIRST(\beta) = \emptyset$
- $A \rightarrow \alpha, A \rightarrow \beta, \varepsilon \in FIRST(\alpha) \implies FIRST(\beta) \cap FOLLOW(A) = \emptyset$

Предположим, что грамматика LL1, но одно из предположений не выполняется.

Пусть первое неверно:

- $FIRST(\alpha) \cap FIRST(\beta) \ni c \in \Sigma$
 $S \xRightarrow{*} xA\xi \xRightarrow{*} x\alpha\xi \xRightarrow{*} xc\alpha'\xi$
 $S \xRightarrow{*} xA\xi \xRightarrow{*} x\beta\xi \xRightarrow{*} xc\beta'\xi$
- $\varepsilon \in FIRST(\alpha) \cap FIRST(\beta)$
 $S \xRightarrow{*} xA\xi \xRightarrow{*} x\alpha\xi \xRightarrow{*} x\xi \xRightarrow{*} xc\eta$
 $S \xRightarrow{*} xA\xi \xRightarrow{*} x\beta\xi \xRightarrow{*} x\xi \Rightarrow xc\eta$

Пусть второе неверно:

$\varepsilon \in FIRST(\alpha) \quad c \in FIRST(\beta)$

...

Пусть $\Gamma \notin LL(1)$

...

Литература:

Ахо, Сети, Ульман "Компиляторы: принципы построения, ..." на ней нарисован дракон (Dragon book)

Две проблемы, которые делают грамматику не LL1 (есть ещё, но вот есть вот эти две):

- Если в грамматике Γ есть левая рекурсия

$$A \xRightarrow{+} A\alpha$$
- Если в грамматике Γ есть правое ветвление

$$A \rightarrow \alpha\beta \quad A \rightarrow \alpha\gamma \quad \exists x \neq \varepsilon : \alpha \xRightarrow{*} x$$

2

✓ Методы Трансляции с 2022-09-13

3

...

$A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$

```
FIRST_1(alpha) = FIRST(alpha) \ eps U (FOLLOW(A) if eps \in FIRST(alpha))
```

```
Node A()
  Node res = Node('A')
  switch (token)
    case FIRST_1(alpha_1)
      // alpha_1[1] \in N, B = alpha_1[1]
      Node n1 = B()
      res.children.push_back(n1)
      // alpha_1[2] \in \Sigma, c = alpha_1[2]
      ensure(token = c; "expected".c."found".token)
      Node n2 = Node(token)
      res.children.push_back(n2)
      next_token()
    ... // end of switch
  default:
    error("unexpected".token)
  return res
```

```
E  → TE'
E' → +TE'
E' → eps
T  → FT'
T' → * FT'
T' → eps
```

F → n
F → (E)

```
Node E()
Node res = Node('E')
switch(token):
    case '(', 'n':
        n1 = T()
        res.children.push_back(n1)
        n2 = E'()
        res.children.push_back(n2)
        break
    default:
        error( ... )
return res

Node E'()
Node res = Node('E')
switch(token)
    case '+':
        ensure (token = '+')
        res.children.push_back(Node('+'))
        next_token()
        res.children.push_back(T())
        res.children.push_back(E'())
        break
    case '$', ')':
        break
    default:
        error
return res

Node F()
Node res = Node(F)
switch(token)
    case 'n':
        res.children.push_back(n)
        nextToken()
        break
    case '(':
        res.children.push_back(Node('('))
        nextToken()
        res.children.push_back(E)
        ensure(token = ')')
        res.children.push_back(Node(')'))
        nextToken()
```

A → beta A'
A' → alpha A'

$A' \rightarrow \epsilon$

```
Node A()
    Node res = Node('A')
    process beta
    while (token \in FIRST(alpha))
        t = res
        res = Node('A')
        res.children.push_back(t)
    process alpha
    ensure (token \in FOLLOW(A))
    return res
```

3

☑ Методы Трансляции с 2022-09-20