

# Формальні граматики

Матеріал з Вікіпедії — вільної енциклопедії.

**Формальна граMATика** або просто **граMATика** в теорії формальних мов — спосіб опису формальної мови, тобто виділення деякої підмножини з множини всіх слів деякого скінченного алфавіту. Розрізняють *породжувальні* і *аналітичні* граматики — перші ставлять правила, за допомогою яких можна побудувати будь-яке слово мови, а другі дозволяють по даному слову визначити, входить воно в мову чи ні. Формальні граматики були введені американським вченим, математиком та філософом, Н. Хомським у 50-тих роках XX сторіччя.

## Зміст

Поняття алфавіту

Формальна граMATика

Виведення в формальних граMATиках

Неоднозначності, та стратегії виводу

Формальні мови

Класифікація граMATик

Класифікація за Хомським

Співвідношення між типами граMATик

Співвідношення між типами мов

Використання формальних граMATик

Див. також

Посилання

## Поняття алфавіту

**Алфавіт** — скінченна множина символів. *ε* — порожній ланцюжок, слово, послідовність. Алфавітом є об'єднання алфавітів, перетин, різниця алфавітів. Нехай *T* — алфавіт, тоді:

- T*<sup>+</sup> — множина усіх можливих послідовностей, що складені з елементів цього алфавіту крім порожньої послідовності *ε*.
- T*<sup>\*</sup> — множина усіх можливих послідовностей, що складені з елементів цього алфавіту, будь-якої довжини. Отже: *T*<sup>\*</sup> = *T*<sup>+</sup> ∪ {*ε*}
- T*<sup>k</sup> — множина усіх можливих послідовностей, що складені з елементів цього алфавіту, довжини не більше *k*.

**Мова** в алфавіті *T* це множина ланцюжків скінченної довжини в цьому алфавіті. Зрозуміло, що кожна мова в алфавіті *T* є підмножиною множини *T*<sup>\*</sup>.

## Формальна граMATика

**Формальна граMATика** - це четвірка ***G* = {*N*, *T*, *P*, *S*}**. Де:

- T*** — алфавіт термінальних символів, терміналів (від англ. terminate - завершитись). Термінальні символи є алфавітом мови.
- N*** — алфавіт нетермінальних символів, нетерміналів. *T* ∩ *N* = ∅; Нетермінали не входять в алфавіт мови.
- S*** — аксіома, спеціально виділений нетермінальний символ з якого починається опис граматики.

*S* ∈ *N*

- $P$  — скінченна підмножина множини  $(T \cup N)^+ \times (T \cup N)^*$ . Інколи визначають так:  
 $\alpha \in (T \cup N)^* \times N \times (T \cup N)^*, \beta \in (T \cup N)^*$ .

Елемент  $(\alpha, \beta)$  з множини  $P$  називається *правилом виводу* і записується у вигляді  $\alpha \rightarrow \beta$ . Таким чином, ліва частина правила не може бути порожньою. Правила з однаковою лівою частиною записують:  $\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n, \beta_i, i = 1, 2, \dots, n$  - називаються альтернативами правила виводу з ланцюжка  $\alpha$ .

## Виведення в формальних граматиках

Ланцюжок  $\beta \in (T \cup N)^*$  назвемо **безпосередньо виведеним** з ланцюжка  $\alpha \in (T \cup N)^+$  в граматиці  $G = \{T, N, S, P\}$  (позначається  $\alpha \Rightarrow \beta$ ), якщо  $\alpha = \xi_1 \gamma \xi_2, \beta = \xi_1 \delta \xi_2 : \xi_1, \xi_2, \delta \in (T \cup N)^*, \gamma \in (T \cup N)^+, \gamma \rightarrow \delta \in P$ .

Ланцюжок  $\beta \in (T \cup N)^*$  назвемо **виведеним** з ланцюжка  $\alpha \in (T \cup N)^+$  в граматиці  $G = \{T, N, S, P\}$  (позначається  $\alpha \Rightarrow^* \beta$ ), якщо

$$\exists \gamma_0, \gamma_1, \dots, \gamma_n, n \geq 0 : \alpha = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \dots \Rightarrow \gamma_n = \beta$$

Термінальний рядок називається **правильним** (або *сентенціальною формою*) відносно граматики  $G$  якщо він виводиться з аксіоми цієї граматики.

## Неоднозначності, та стратегії виводу

ГраMATика  $G$  називається **неоднозначною**, якщо існує декілька варіантів виводу слова  $\omega$  в  $G$  ( $\omega \in L(G)$ ).

**Приклад.** Розглянемо таку граматику  $G = \langle N, \Sigma, P, S \rangle$  зі схемою  $P$ .

$S \rightarrow S + S \mid S * S \mid a$

Покажемо, що для ланцюжка  $\omega = a + a + a$  існує щонайменше **два** варіанти виводу:

$$1. S \Rightarrow S + S \Rightarrow S + S + S \Rightarrow a + S + S \Rightarrow a + a + S \Rightarrow a + a + a$$

$$2. S \Rightarrow S + S \Rightarrow a + S \Rightarrow a + S + S \Rightarrow a + a + S \Rightarrow a + a + a$$

В теорії граматики розглядається декілька стратегій виведення ланцюжка  $\omega$  в  $G$ .

**Лівостороння стратегія виводу** ланцюжка  $\omega$  в  $G$  — це послідовність кроків безпосереднього виводу, при якій на кожному кроку до уваги береться перший зліва направо нетермінал.

**Правостороння стратегія виводу**  $\omega$  в  $G$  протилежна лівосторонній стратегії. З виводом  $\omega$  в  $G$  пов'язане синтаксичне дерево, яке визначає синтаксичну структуру програми.

## Формальні мови

Формальна мова породжена граматиною  $G$  - це множина термінальних рядків, що виводяться з аксіоми, тобто множина усіх правильних рядків.

$$L(G) = \{\alpha \in T^* \mid S \Rightarrow^* \alpha\}$$

З іншого боку, множина термінальних рядків, таких, що вони можуть бути виведені в граматиці  $G$ , називається мовою, що може бути розпізнана в даній граматиці, або допускається даною граматиною.

- Граматики  $G_1$  та  $G_2$  називаються **еквівалентними**, якщо  $L(G_1) = L(G_2)$ .
- Граматики  $G_1$  та  $G_2$  називаються **майже еквівалентними**, якщо  $L(G_1) \cup \{\epsilon\} = L(G_2) \cup \{\epsilon\}$ , тобто мови, ними породжувані, відрізняються не більш ніж на  $\epsilon$ .

ГраMATика може бути **породжуючою** та **розпізнавальною**, це залежить від "напрямку" застосування правил. Для породжуючих грамаТИК виведення починається з аксіоми і закінчується термінальним рядком (рядком, що складається тільки з терміналів). А розпізнавальні аналізують вхідний термінальний рядок на правильність, чиможна такий рядок вивести в цій грамаТИЦІ. Коротко кажучи - породжуючі це "згори вниз", а розпізнавальні - "знизу вгору". Остання задача є практичнішою і гострою.

## Класифікація грамаТИК

---

Хомський також запропонував класифікацію грамаТИК. Він виділив чотири типи грамаТИК.

### Класифікація за Хомським

#### Тип 0. ГрамаТИКИ Типу 0

це найзагальніший клас грамаТИК. На правила не накладаються ніяких обмежень, окрім зазначених у визначенні. Вони еквівалентні Машині Тюринга. Доведено існування мов, які породжуються грамаТИКАМИ типу 0, але не грамаТИКАМИ типу 1, але невідомі прості приклади таких мов.

#### Тип 1. ГрамаТИКИ Типу 1

нескорочуюча або контекстно-залежна(КЗ) грамаТИКА. Вибір означення не впливає на множину мов, породжуваних грамаТИКАМИ цього класу, оскільки доведено, що множина мов, породжуваних грамаТИКАМИ, що не укорочують, збігається із множиною мов, породжуваних КЗ-грамаТИКАМИ.

- *Нескорочуючі грамаТИКИ*. Всі правила мають вигляд

$$\begin{aligned} \alpha &\rightarrow \beta \\ \alpha &\in (T \cup N)^+, \beta \in (T \cup N)^* \\ |\alpha| &\leq |\beta| \\ \text{де } |\alpha| &\text{ означає кількість символів у рядку } \alpha. \end{aligned}$$

- *Контекстно-залежні грамаТИКИ*. Всі правила мають вигляд:

$$\begin{aligned} \alpha &\rightarrow \beta \\ \alpha &= \xi_1 A \xi_2, \beta = \xi_1 \gamma \xi_2 : A \in N, \gamma \in (T \cup N)^+, \xi_1, \xi_2 \in (T \cup N)^* \end{aligned}$$

#### Тип 2. ГрамаТИКИ Типу 2. Контекстно-вільні грамаТИКИ(КВ)

Всі правила мають вигляд

$$\begin{aligned} \alpha &\rightarrow \beta \\ \alpha &\in N, \beta \in (T \cup N)^* \end{aligned}$$

Тобто в усіх правилах цього виду зліва стоїть тільки один нетермінал. Тому вони і контекстно вільні, бо на використання правила для даного нетерміналу не впливають символи, що оточують його. Ці символи називають **контекстом**.

#### Тип 3. ГрамаТИКИ Типу 3. Регулярні грамаТИКИ. А-грамаТИКИ

можна визначити як праволінійну або ліволінійну грамаТИКУ, або як змішану. Також ці мови називають скінченно-автоматними, бо вони еквівалентні скінченним автоматом, тобто клас мов, що породжуються грамаТИКАМИ типу 3, збігається з класом мов, які розпізнаються скінченими автоматами. Також ці грамаТИКИ є основою регулярних виразів.

- *Праволінійна грамаТИКА*. Всі правила мають вигляд:

$$\begin{aligned} \alpha &\rightarrow \beta \\ \alpha &\rightarrow \omega\beta, \beta \in N, \omega \in T^* \\ \text{або} \end{aligned}$$

$$\alpha \rightarrow \omega, \omega \in T^*$$

- Ліволінійна граMATика. Всі правила мають вигляд:

$$\alpha \rightarrow \beta$$

$$\alpha \rightarrow \beta\omega, \beta \in N, \omega \in T^*$$

або

$$\alpha \rightarrow \omega, \omega \in T^*$$

Доведено, що праволінійні та ліволінійні граMATики еквівалентні, і існує алгоритм для переведення правил граMATики одного типу в інший тип.

## Співвідношення між типами граMATик

1. Будь-які граMATики типу 3 є граMATиками типу 2.
2. Будь-які КВ-граMATики є граMATиками типу 1 (КЗ або неукорочуючі), з точністю до  $\epsilon$ .
3. Будь-які граMATики типу 1 є граMATиками типу 0.

Мова  $L(G)$  є мовою типу  $k$ , якщо її можна описати граMATикою типу  $k$ .

## Співвідношення між типами мов

1. кожна регулярна мова є КВ-мовою, але існують КВ- мови, що не є регулярними (наприклад,  $L = \{a^n b^n \mid n > 0\}$ );
2. кожна КВ- мова є КЗ- мовою, але існують КЗ-мови, що не є КВ-мовами ( наприклад,  $L = \{a^n b^n c^n \mid n > 0\}$ );
3. кожна КЗ-мова є мовою типу 0.

Варто підкреслити, що якщо мова задана граMATикою типу  $k$ , то це не означає, що не існує граMATики типу  $k'$  ( $k' > k$ ), що опише ту ж мову. Тому, коли говорять про мову типу  $k$ , звичайно мають на увазі максимально можливе  $k$ .

## Використання формальних граMATик

Формальні граMATики це дуже потужний математичний інструмент, який використовується в математичній та комп'ютерній лінгвістиці, описі мов програмування, розробці компіляторів в теорії програмування. Найбільш вживаними є КВ-граMATики і регулярні, бо їх найлегше описати алгоритмічно.

Сама по собі концепція формальних граMATик доволі гнучка, тому не дивно, що з'явилося багато інших інструментів, що розширюють використання та потужність КВ-граMATик і граMATик третього типу. Наприклад, атрибутні граMATики, LL-k граMATики, скінченні автомати, регулярні вирази та множини.

## Див. також

- Нотація Бекуса-Наура
- Синтаксичний аналіз
- Ієрархія Чомські

## Посилання

- Серебряков В.А., Галочкин М.П., Гончар Д.Р., Фуругян М.Г. Теория и реализация языков программирования // М.: МЗ-Пресс, 2006 г, 2-е изд. ISBN 94073-094-9. (рос.)
- Подборка ссылок на литературу по направлению. (рос.)

---

**Цю сторінку востаннє відредаговано о 23:30, 14 грудня 2017.**

Текст доступний на умовах ліцензії [Creative Commons Attribution-ShareAlike](#) також можуть діяти додаткові умови. Детальніше див. [Умови використання](#).