

# Théorie des Langages

## Cours 6 : Grammaires

Lionel Rieg

Grenoble INP - Ensimag, 1<sup>re</sup> année

# Reconnaissance et génération

- **Reconnaissance** : on part de  $w$  et on veut savoir si  $w \in L$ 
  - ▶ automate fini
  - ▶ automate à pile
  - ▶ machine de Turing

⇒ peu de structure pour comprendre le langage

- **Génération** : on veut engendrer un mot d'un langage
  - ▶ expressions régulières
  - ▶ **grammaires**

⇒ règles structurées qui décrivent le langage

Permet de mieux comprendre comment est fait un langage

**Les grammaires sont utilisées pour décrire les langages de programmation.**

## Extrait de la grammaire de Python 3

```
statement: compound_stmt | simple_stmts
simple_stmts:
    | ';' simple_stmt+ [';'] NEWLINE
simple_stmt:
    | assignment
    | return_stmt
    | import_stmt
    | raise_stmt
    | 'pass'
    | assert_stmt
    | 'break'
    | 'continue'
compound_stmt:
    | function_def
    | if_stmt
    | for_stmt
    | try_stmt
    | while_stmt
```

# Intérêt des grammaires

## Un formalisme de description de langage

- permet de définir/caractériser des classes de langages
- algorithmes génériques sur chaque classe
- étudier des propriétés pour chaque classe

## Hiérarchie de Chomsky = compromis expressivité/décidabilité

- Expressivité = quels genres de langages peut-on exprimer ?
- Décidabilité = existe-il un algorithme qui répond au problème ?

## Langage de programmation idéal :

- expressif : sa syntaxe permet beaucoup de chose
- décidable : on doit pouvoir analyser sa syntaxe automatiquement

⇒ **Compromis** : trouver une classe décidable mais suffisamment expressive

# Exemple du français

## Exemple (Grammaire du français)

Une phrase est formée d'un sujet, d'un verbe et, optionnellement, d'un complément.

Un sujet peut être un nom propre ou un nom commun, optionnellement précédé ou suivi par un adjectif.

Un complément est un COD, un COI ou les deux.

Un nom commun peut être « voiture », « maison » ou etc. (tout le vocabulaire).

...

Une grammaire de langue naturelle :

- Explique comment construire des catégories grammaticales ...  
(phrase, sujet, complément, etc.) **vocabulaire non terminal**
- ... en utilisant des **règles**.
- On applique des règles jusqu'à n'avoir plus que des mots  
(voiture, maison, est, etc.). **vocabulaire terminal**

# Définition d'une grammaire

## Définition (grammaire)

Une grammaire est un quadruplet  $G \stackrel{\text{def}}{=} \langle V_T, V_N, S, R \rangle$  où

- $V_T$  est un ensemble fini appelé le **vocabulaire terminal**
- $V_N$  est un ensemble fini appelé le **vocabulaire non terminal**
- $S \in V_N$  est l'**axiome**
- $R$  est un ensemble de **règles** de la forme  $u \rightarrow v$ , avec
  - ▶  $u \in V^+$
  - ▶  $v \in V^*$

On pose  $V \stackrel{\text{def}}{=} V_N \cup V_T$  et on demande  $V_N \cap V_T = \emptyset$ .

Comment s'utilise une grammaire ?

- on commence avec l'axiome  $S$  (Start symbol)
- on applique des règles tant qu'on veut
- on peut s'arrêter si le mot ne contient que des symboles terminaux (mais on peut aussi continuer)

# Notations

## Conventions

- On distingue graphiquement  $V_T$  et  $V_N$  (**gras**, souligné, *italique*, etc.)
- $V$  est l'ensemble des symboles présents dans  $G$  et  $V_N = V \setminus V_T$ .
- L'axiome  $S$  est la partie gauche de la première règle.
- Si plusieurs règles  $u \rightarrow w_1, \dots, u \rightarrow w_n$ ,  
on peut les noter  $u \rightarrow w_1 \mid \dots \mid w_n$ .
- Il arrive qu'on ne précise pas  $V_T$ , s'il se déduit du langage.

## Notations

- $\rightarrow$  pour une règle,  $\Rightarrow$  pour son application dans un mot
- Pour mieux identifier l'utilisation d'une règle dans un mot,  
on souligne parfois sa partie gauche :  $ab\underline{S}c \Rightarrow abcac$

## Exemple

### Exemple (Grammaire pour $(ab)^*(ba)^+$ )

$$V_T = \{a, b\} \quad V_N = \{X, AB, BA\} \quad S = X$$

$$R = \left\{ X \xrightarrow{1} AB BA, \quad AB \xrightarrow{2} abAB, \quad AB \xrightarrow{3} \varepsilon, \quad BA \xrightarrow{4} baBA, \quad BA \xrightarrow{5} ba \right\}$$

Quelques dérivations possibles :

- $\underline{X} \xrightarrow{1} \underline{AB} BA \xrightarrow{3} \varepsilon \underline{BA} \xrightarrow{5} \varepsilon ba = ba$
- $\underline{X} \xrightarrow{1} \underline{AB} BA \xrightarrow{2} ab \underline{AB} BA \xrightarrow{2} abab \underline{AB} BA$   
 $\xrightarrow{3} abab \underline{BA} \xrightarrow{4} ababba \underline{BA} \xrightarrow{5} ababbaba$
- $\underline{X} \xrightarrow{1} AB \underline{BA} \xrightarrow{4} \underline{AB} ba BA \xrightarrow{2} ab \underline{AB} ba BA$   
 $\xrightarrow{2} abab AB ba \underline{BA} \xrightarrow{5} abab \underline{AB} baba \xrightarrow{3} ababbaba$



# Premiers exercices

## Exercice

Donner des grammaires pour les langages suivants :

- $\{a^n b^n \mid n \in \mathbb{N}\}$
- $\{a^n b^n c^n \mid n \geq 1\}$
- $\{a^n b^n \mid n \in \mathbb{N}\}$  :
  - ▶  $V_T \stackrel{\text{def}}{=} \{a, b\}$
  - ▶  $V_N \stackrel{\text{def}}{=} \{S\}$
  - ▶  $S \stackrel{\text{def}}{=} S$
  - ▶  $R \stackrel{\text{def}}{=} \{S \rightarrow aSb, S \rightarrow \varepsilon\} = \{S \rightarrow aSb \mid \varepsilon\}$
- $\{a^n b^n c^n \mid n \geq 1\}$  :
  - ▶  $V_T = \{a, b, c\}$
  - ▶  $S \rightarrow aSBc \mid abc$
  - ▶  $bB \rightarrow bb$
  - ▶  $cB \rightarrow Bc$
  - ▶ donc par convention :  $V_N = \{S, B\}$  et l'axiome est S

# Relation de dérivation

## Définition (Dérivation)

Soit  $G = \langle V_T, V_N, S, R \rangle$  une grammaire.

Étant donné  $x, y \in V^*$ , on dit que  $x$  **dérive vers**  $y$ , noté  $x \Longrightarrow y$ , lorsque :

- il existe une règle  $u \rightarrow v$  dans  $R$ ,
- $x$  s'écrit  $\alpha u \beta$  pour  $\alpha, \beta \in V^*$
- $y$  vaut  $\alpha v \beta$ .

Idée : **on a remplacé  $u$  par  $v$  dans  $x$  pour obtenir  $y$ .**

Pour être plus précis, on peut préciser la règle utilisée et sa position :

$$\alpha \underline{u} \beta \xrightarrow{u \rightarrow v} \alpha v \beta$$



Ne pas confondre  $\rightarrow$  (règle) et  $\Longrightarrow$  (dérivation).

Ne pas confondre  $\Longrightarrow$  (dérivation) et  $\implies$  (implication logique).

# Relation de dérivation itérée

## Définition

- On note  $\Longrightarrow^p$  les dérivations de longueur  $p$ , définies récursivement sur  $p$  par :
  - $u \Longrightarrow^0 v \stackrel{\text{def}}{=} u = v$
  - $u \Longrightarrow^{p+1} v \stackrel{\text{def}}{=} \exists w, u \Longrightarrow w \wedge w \Longrightarrow^p v$
- $\Longrightarrow^*$  est la fermeture réflexive transitive de  $\Longrightarrow$  :  $\Longrightarrow^* \stackrel{\text{def}}{=} \bigcup_{p \in \mathbb{N}} \Longrightarrow^p$

## Composition des dérivations :

Si  $u_1 \Longrightarrow^{p_1} v_1$  et  $u_2 \Longrightarrow^{p_2} v_2$ , alors  $u_1 u_2 \Longrightarrow^{p_1+p_2} v_1 v_2$



Ce n'est souvent pas la seule dérivation possible dans  $u_1 u_2$  !

**Exercice :** Donner une définition inductive de  $u \Longrightarrow^* v$ .

**Base :**  $u = v$

**Induction :**  $\exists w, u \Longrightarrow w \Longrightarrow^* v$

# Langage engendré

## Définition (Langage engendré)

Soit  $G = \langle V, T, V_N, S, R \rangle$  une grammaire.

Pour  $u \in V^* = (V_T \cup V_N)^*$ , on définit le **langage engendré** par  $u$  :

$$\mathcal{L}(u) \stackrel{\text{def}}{=} \{w \in V_T^* \mid u \Longrightarrow^* w\}$$

Le **langage d'une grammaire** est le langage engendré par son axiome :

$$\mathcal{L}(G) \stackrel{\text{def}}{=} \mathcal{L}(S) = \{w \in V_T^* \mid S \Longrightarrow^* w\}$$

## Définition (Grammaires équivalentes)

Comme pour les automates, deux grammaires  $G_1$  et  $G_2$  sont dites **équivalentes** si elles engendrent le même langage :  $\mathcal{L}(G_1) = \mathcal{L}(G_2)$ .

## Exemples de langage engendré

- Considérons la grammaire suivante :
$$\begin{aligned}S &\rightarrow AB \\A &\rightarrow a \mid aA \\B &\rightarrow bB \mid \varepsilon\end{aligned}$$

Que valent  $V_T$ ,  $V_N$ ,  $V$ ,  $S$  ?

$$V_T = \{a, b\} \quad V_N = \{S, A, B\} \quad V = \{a, b, S, A, B\}$$

Quels sont les langages engendrés par  $A$ ,  $B$ ,  $S$ ,  $AbA$  ?

$$\mathcal{L}(A) = a^+ \quad \mathcal{L}(B) = b^* \quad \mathcal{L}(S) = a^+b^* \quad \mathcal{L}(AbA) = a^+ba^+$$

- Autre grammaire :
$$\begin{aligned}S &\rightarrow aSBc \mid \varepsilon \\Bc &\rightarrow bBc \mid bc \mid b \\cb &\rightarrow bc\end{aligned}$$

$$V_T = \{a, b, c\} \quad V_N = \{S, B\}$$

$$\mathcal{L}(B) = \emptyset \quad \mathcal{L}(b) = \{b\} \quad \mathcal{L}(cb) = \{cb, bc\}$$

$$\mathcal{L}(Bc) = b^+(c + \varepsilon)$$

$$\mathcal{L}(S) = \{a^n w \mid n \in \mathbb{N}, |w|_c \leq n \leq |w|_b, \forall x \text{ préfixe de } w, |x|_b \geq |w|_c\}$$

# Énumération

Une grammaire permet d'énumérer des mots. On peut définir l'ensemble  $\text{Gen}_p$  des mots sur  $V$  engendrés après  $p$  dérivations par récurrence sur  $p$  :

- On commence avec le singleton contenant l'axiome :  $\text{Gen}_0 = \{S\}$
- On applique toutes les règles possibles à toutes les positions possibles sur les mots de  $\text{Gen}_p$  :  $\text{Gen}_{p+1} = \{w \in V \mid \exists u \in \text{Gen}_p, u \Longrightarrow w\}$

## Exemple

$S \rightarrow AB \quad A \rightarrow a \mid aA \quad B \rightarrow bB \mid \varepsilon$

- $\text{Gen}_0 = \{S\}$
- $\text{Gen}_1 = \{AB\}$
- $\text{Gen}_2 = \{aB, aAB, AbB, A\}$
- $\text{Gen}_3 = \{abB, a, aaB, aaAB, aAbB, aA, abB, aAbB, AbbB, Ab, a, aA\}$

$\Rightarrow$  Procédure de **semi-décision** pour le problème  $w \in \mathcal{L}(G)$  ?  
peu efficace, semi-decision car on ne sait pas quand s'arrêter

Lien avec le théorème du point fixe :  $\bigcup_{p \leq n} \text{Gen}_p = f^{n+1}(\emptyset)$   
avec  $f : X \mapsto \{w \mid u \in X \wedge u \Longrightarrow w\} \cup \{S\}$

# Hiérarchie de Chomsky

Classification des grammaires par **restriction sur les règles**

- **type 3** : grammaire régulière la plus décidable
- **type 2** : grammaire hors-contexte (HC)
- **type 1** : grammaire sous-contexte
- **type 0** : grammaire générale la plus expressive

Étude du **compromis expressivité/décidabilité**

Propriétés étudiées :

- $w \in \mathcal{L}(G)$  ?
- $\mathcal{L}(G) = \emptyset$  ?
- $|\mathcal{L}(G)| = \infty$  ?
- $\mathcal{L}(G_1) = \mathcal{L}(G_2)$  ?

## Grammaire régulière (type 3)

- **Restriction** : règle de la forme

$$A \rightarrow xB \quad \text{ou} \quad A \rightarrow \varepsilon, \quad \text{avec} \quad A, B \in V_N, x \in V_T$$

### Grammaire **linéaire à droite**

- ▶ linéaire : un seul non-terminal à droite de la règle
- ▶ à droite : ce non-terminal est en dernière position  
(à gauche = en première position)

- **Reconnaisseurs** : automates finis
- **Idée** : non-terminal = état de l'automate (sans  $\varepsilon$ -transition)
- **Exemple** :  $A \rightarrow aA \mid bB \quad B \rightarrow bB \mid \varepsilon$   
 $\mathcal{L}(A) = a^*b^+$
- **Caractérisation équivalente** :

$$A \rightarrow wB \quad \text{ou} \quad A \rightarrow w, \quad \text{avec} \quad A, B \in V_N, w \in V_T^*$$



# Conversion grammaire régulière / automate fini



L'automate ne doit avoir qu'un unique état initial et pas d' $\varepsilon$ -transition (mais on sait faire).

## Automate fini

Alphabet  $V$

États  $Q$

Transition  $(q, x, q') \in \delta \iff$

État initial  $i$

États acceptants  $q \in F$

## Grammaire régulière

Symboles terminaux  $V_T$

Symboles non-terminaux  $V_N$

Règle  $q \rightarrow xq'$

Axiome  $S$

Règle  $q \rightarrow \varepsilon$

Invariant de correction : (preuve par induction)

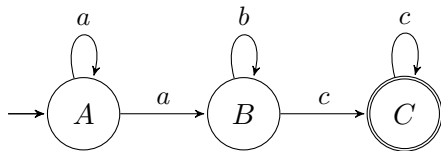
$\chi$  chemin de  $q$  à  $q'$  de trace  $w \iff q \implies^* wq'$

## Exercice : faire les conversions suivantes

- Automate fini vers grammaire régulière :  $a^+b^*c^+$
- Grammaire régulière vers automate fini :  $P \rightarrow aI \mid bP$   
 $I \rightarrow aP \mid bI \mid \varepsilon$

## Correction de l'exercice

- Automate fini vers grammaire régulière :  $a^+b^*c^+$



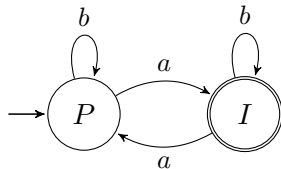
$\sim$

$A \rightarrow$	$aA$	$ $	$aB$
$B \rightarrow$	$bB$	$ $	$cC$
$C \rightarrow$	$cC$	$ $	$\varepsilon$

- Grammaire régulière vers automate fini :

$P \rightarrow$	$aI$	$ $	$bP$
$I \rightarrow$	$aP$	$ $	$bI$
			$\varepsilon$

$\sim$



## Grammaire hors-contexte/non-contextuelle (type 2)

- **Restriction** : règle de la forme

$$A \rightarrow w \quad \text{avec} \quad A \in V_N, w \in V^*$$

- **Idée** : pas de contexte autour du non-terminal à transformer
- **Exemple** :  $S \rightarrow aSb \mid \varepsilon$   
 $\mathcal{L}(S) = \{a^n b^n \mid n \in \mathbb{N}\}$
- **Reconnaisseurs** : automates à pile

**Bon compromis expressivité/décidabilité**

⇒ algorithmes efficaces de reconnaissance

⇒ utilisées pour les langages de programmation

La classe qu'on étudiera le plus.

# Grammaire sous-contexte/contextuelle (type 1)

- **Restriction** : règle de la forme

$$\alpha A \beta \rightarrow \alpha w \beta \quad \text{avec} \quad A \in V_N, \quad \alpha, \beta, w \in V^* \quad \text{et} \quad w \neq \varepsilon$$

- **Idée** : non-terminal à transformer, avec un contexte
- **Caractérisation équivalente** (plus utile) :

$$u \rightarrow v \quad \text{avec} \quad u, v \in V^+ \quad \text{et} \quad |u| \leq |v|$$

- **Exemple** :  $S \rightarrow aSBc \mid abc \quad cB \rightarrow Bc \quad aB \rightarrow ab$   
 $\mathcal{L}(S) = \{a^n b^n c^n \mid n \geq 1\}$
- **Reconnaisseurs** : machines de Turing linéairement bornée

# Propriétés des grammaires sous-contexte

- La caractérisation  $|u| \leq |v|$  implique que la taille d'un mot ne peut pas décroître au cours d'une dérivation.  
 $\Rightarrow w \in \mathcal{L}(S)$  devient décidable : on peut s'arrêter lorsque la taille des mots engendrés dépasse celle du mot à tester.
  1. Commencer avec  $\text{Gen}_0 = \{S\}$
  2. Calculer  $D_p = \{x \in \text{Gen}_p \mid |x| \leq |w|\}$
  3. Jusqu'à trouver  $w$  ou qu'il soit vide.
- Lorsque  $\varepsilon$  doit être dans le langage d'une grammaire sous-contexte, on ajoute :
  - ▶  $Z \in V_N$ , nouvel axiome
  - ▶  $Z \rightarrow S \mid \varepsilon$  avec  $S$  l'ancien axiome

# Grammaire générale (type 0)

- Restriction : aucune
- Reconnaisseurs : machine de Turing

Les plus expressives, mais peu de chose décidables

$\Rightarrow$  même  $w \in \mathcal{L}(G)$  n'est plus décidable

# Relations entre classes

## Inclusions :

- régulière  $\subset$  hors-contexte
- hors-contexte sans  $\varepsilon$ -règle  $\subset$  sous-contexte ( $\varepsilon$ -règle =  $u \rightarrow \varepsilon$ )
- sous-contexte  $\subset$  générale

## Classes de langages

Un langage est dit **régulier/hors-contexte/sous-contexte** s'il peut être généré par une grammaire de même type.

Pas d'algorithme pour décider de la classe d'un langage

## Décidabilité de propriétés :

Propriétés \ type de $G$	3 (reg)	2 (HC)	1 (SC)	0 (gen)
$w \in \mathcal{L}(G) ?$	✓	✓	✓	✗
$\mathcal{L}(G) = \emptyset ?$	✓	✓	✗	✗
$ \mathcal{L}(G)  = \infty ?$	✓	✓	✗	✗
$\mathcal{L}(G_1) = \mathcal{L}(G_2) ?$	✓	✗	✗	✗