

# Languages et expressions régulières

---

Ambre Le Berre

2025/2026

MP option info

# Agenda

1. Vocabulaire
2. Expression régulière, langages réguliers

# Vocabulaire

---

# Alphabet et mots

- Un **alphabet** est un ensemble  $\Sigma$  de **lettres**. Par exemple,  $\Sigma = \{a, b\}$  est un alphabet.
- 
- 
- 
-

# Alphabet et mots

- Un **alphabet** est un ensemble  $\Sigma$  de **lettres**. Par exemple,  $\Sigma = \{a, b\}$  est un alphabet.
- Un **mot** sur un alphabet est une suite finie de lettres de cet alphabet. Par exemple, *abba* est un mot sur  $\Sigma$ .
- 
- 
-

# Alphabet et mots

- Un **alphabet** est un ensemble  $\Sigma$  de **lettres**. Par exemple,  $\Sigma = \{a, b\}$  est un alphabet.
- Un **mot** sur un alphabet est une suite finie de lettres de cet alphabet. Par exemple, *abba* est un mot sur  $\Sigma$ .
- On note  $\varepsilon$  le mot vide, qui ne contient aucune lettre.
- 
-

# Alphabet et mots

- Un **alphabet** est un ensemble  $\Sigma$  de **lettres**. Par exemple,  $\Sigma = \{a, b\}$  est un alphabet.
- Un **mot** sur un alphabet est une suite finie de lettres de cet alphabet. Par exemple, *abba* est un mot sur  $\Sigma$ .
- On note  $\varepsilon$  le mot vide, qui ne contient aucune lettre.
- Pour deux mots  $u$  et  $v$ , on note  $uv$  leur concaténation.
- On note  $|u| \in \mathbb{N}$  la longueur d'un mot.

# Vocabulaires sur les mots

Soit  $m = m_1 m_2 \dots m_n \in \Sigma^*$  un mot.

- un **préfixe** de  $m$  est un mot  $m_1 \dots m_i$  avec  $0 \leq i \leq n$ .

- 

- 

- 

-

# Vocabulaires sur les mots

Soit  $m = m_1 m_2 \dots m_n \in \Sigma^*$  un mot.

- un **préfixe** de  $m$  est un mot  $m_1 \dots m_i$  avec  $0 \leq i \leq n$ .
- un **suffixe** de  $m$  est un mot  $m_{i+1} \dots m_n$  avec  $0 \leq i \leq n$ .
- 
- 
-

# Vocabulaires sur les mots

Soit  $m = m_1 m_2 \dots m_n \in \Sigma^*$  un mot.

- un **préfixe** de  $m$  est un mot  $m_1 \dots m_i$  avec  $0 \leq i \leq n$ .
- un **suffixe** de  $m$  est un mot  $m_{i+1} \dots m_n$  avec  $0 \leq i \leq n$ .
- un **facteur** de  $m$  est un mot  $m_i \dots m_j$  avec  $0 \leq i, j \leq n$ .
- 
-

# Vocabulaires sur les mots

Soit  $m = m_1 m_2 \dots m_n \in \Sigma^*$  un mot.

- un **préfixe** de  $m$  est un mot  $m_1 \dots m_i$  avec  $0 \leq i \leq n$ .
- un **suffixe** de  $m$  est un mot  $m_{i+1} \dots m_n$  avec  $0 \leq i \leq n$ .
- un **facteur** de  $m$  est un mot  $m_i \dots m_j$  avec  $0 \leq i, j \leq n$ .
- un **sous-mot** de  $m$  est un mot  $m_{i_1} m_{i_2} \dots m_{i_k}$  avec  $i_1 < i_2 < \dots < i_k \in \llbracket 1, n \rrbracket$
-

# Vocabulaires sur les mots

Soit  $m = m_1 m_2 \dots m_n \in \Sigma^*$  un mot.

- un **préfixe** de  $m$  est un mot  $m_1 \dots m_i$  avec  $0 \leq i \leq n$ .
- un **suffixe** de  $m$  est un mot  $m_{i+1} \dots m_n$  avec  $0 \leq i \leq n$ .
- un **facteur** de  $m$  est un mot  $m_i \dots m_j$  avec  $0 \leq i, j \leq n$ .
- un **sous-mot** de  $m$  est un mot  $m_{i_1} m_{i_2} \dots m_{i_k}$  avec  $i_1 < i_2 < \dots < i_k \in \llbracket 1, n \rrbracket$
- $\varepsilon$  est préfixe, suffixe, facteur et sous-mot de n'importe quel mot.

# Languages

- Un **language** est un ensemble  $L \subseteq \Sigma^*$  de mots sur un alphabet
-

# Languages

- Un **language** est un ensemble  $L \subseteq \Sigma^*$  de mots sur un alphabet
- On note  $\emptyset$  le language vide (à ne pas confondre avec  $\{\varepsilon\}$  )

# Expression régulière, langages réguliers

---

# Opérations régulières sur les langages

On définit plusieurs **opérations** dites **régulières** sur les langages. Soit  $L_1, L_2$  deux langages sur un alphabet  $\Sigma$ .

- $L_1 \cup L_2 = \{u \mid u \in L_1 \text{ ou } u \in L_2\}$  est l'**union** des langages.

- 

- 

-

# Opérations régulières sur les langages

On définit plusieurs **opérations** dites **régulières** sur les langages. Soit  $L_1, L_2$  deux langages sur un alphabet  $\Sigma$ .

- $L_1 \cup L_2 = \{u \mid u \in L_1 \text{ ou } u \in L_2\}$  est l'**union** des langages.
- $L_1 L_2 = \{uv \mid u \in L_1, v \in L_2\}$  est la **concaténation** des langages.

•

•

# Opérations régulières sur les langages

On définit plusieurs **opérations** dites **régulières** sur les langages. Soit  $L_1, L_2$  deux langages sur un alphabet  $\Sigma$ .

- $L_1 \cup L_2 = \{u \mid u \in L_1 \text{ ou } u \in L_2\}$  est l'**union** des langages.
- $L_1 L_2 = \{uv \mid u \in L_1, v \in L_2\}$  est la **concaténation** des langages.
- $L_1^* = \{u_1 u_2 \dots u_n \mid n \in \mathbb{N}, u_1 \dots u_n \in L_1\}$  l'**étoile de Kleene** d'un langage (ou juste **étoile**).
-

# Opérations régulières sur les langages

On définit plusieurs **opérations** dites **régulières** sur les langages. Soit  $L_1, L_2$  deux langages sur un alphabet  $\Sigma$ .

- $L_1 \cup L_2 = \{u \mid u \in L_1 \text{ ou } u \in L_2\}$  est l'**union** des langages.
- $L_1 L_2 = \{uv \mid u \in L_1, v \in L_2\}$  est la **concaténation** des langages.
- $L_1^* = \{u_1 u_2 \dots u_n \mid n \in \mathbb{N}, u_1 \dots u_n \in L_1\}$  l'**étoile de Kleene** d'un langage (ou juste **étoile**).
- L'intersection et la soustraction ensembliste ne sont pas des opérations régulières !

# Expressions régulières

On utilise **expressions régulières** pour définir une partie des langages.

## Définition inductive

- Cas de bases :
  - $\varepsilon$  est une expression régulière représentant le langage  $\{\varepsilon\}$
  - si  $a \in \Sigma$ ,  $a$  est une expression régulière représentant le langage  $\{a\}$ .
-

# Expressions régulières

On utilise **expressions régulières** pour définir une partie des langages.

## Définition inductive

- Cas de bases :
  - $\varepsilon$  est une expression régulière représentant le langage  $\{\varepsilon\}$
  - si  $a \in \Sigma$ ,  $a$  est une expression régulière représentant le langage  $\{a\}$ .
- Cas inductifs : si  $e_1, e_2$  sont des expressions régulières qui représentent des langages  $L_1$  et  $L_2$ , alors
  - $e_1^*$  représente le langage  $L_1^*$
  - $e_1 \mid e_2$  représente le langage  $L_1 \cup L_2$
  - $e_1 e_2$  représente le langage  $L_1 L_2$

# Expressions régulières : raccourcis

On utilise souvent les syntaxes suivantes par facilité :

- $e^+$  pour désigner  $ee^*$
- $e^n$  pour désigner  $eee\dots e$  (n fois)
- $e?$  pour désigner  $e \mid \varepsilon$

# Quelques exemples

Décrivez, en quelques mots, les langages décrits par les expressions régulières suivantes sur l'alphabet  $\Sigma = \{a, b\}$ :

- $(a|b)^*$
- $(a|b)^+$
- $(b^*ab^*)^5$
- $(ab|b)^*$

# Dans l'autre sens

Donnez des expressions régulières pour les langages suivants :

- les mots qui contiennent “abba”
- les mots qui contiennent d'abord que des b (au moins 1), puis que des a (au moins 1)
- les mots dont la longueur est multiple de 3
- les mots dont “abba” est un sous-mot

# Language régulier / rationnel

Un langage **régulier** ou **rationnel** est un langage qui peut être représenté par une expression régulière.

Q : Est ce que tous les langages sont réguliers ?

# Language régulier / rationnel

Un langage **régulier** ou **rationnel** est un langage qui peut être représenté par une expression régulière.

Q : Est ce que tous les langages sont réguliers ? Non !

# Un exemple de langage non régulier

Le langage  $L_0 = \{a^n b^n \mid n \in \mathbb{N}\}$  n'est pas régulier.

# Un exemple de langage non régulier

Le langage  $L_0 = \{a^n b^n \mid n \in \mathbb{N}\}$  n'est pas régulier.

**Intuition** : les expressions régulières ont toujours une “mémoire” finie. Entre autre, elle ne peuvent pas “compter”  $n$  “a” pour attendre ensuite le même nombre de “b”.

# Lemme de l'étoile (première version)

Soit  $L$  un langage régulier. Alors il existe  $N$ , tel que pour tout mot  $u \in L$  avec  $|u| \geq N$ , il existe une décomposition  $u = xyz$ , telle que :

- $|xy| \leq N$
- $y \neq \varepsilon$
- $xy^*z \subseteq L$

# Montrer qu'un langage n'est pas rationnel

Montrons que  $L_0$  n'est pas rationnel par l'absurde : si  $L_0$  est rationnel, alors, soit  $N \in \mathbb{N}$  suffisamment grand pour appliquer le lemme de l'étoile.

On prend le mot  $u = a^{N+1}b^{N+1} \in L_0$ . Supposons qu'il existe  $xyz = u$  tel que

$$|xy| \leq N, y \neq \varepsilon \text{ et } xy^*z \subseteq L$$

# Montrer qu'un langage n'est pas rationnel

Montrons que  $L_0$  n'est pas rationnel par l'absurde : si  $L_0$  est rationnel, alors, soit  $N \in \mathbb{N}$  suffisamment grand pour appliquer le lemme de l'étoile.

On prend le mot  $u = a^{N+1}b^{N+1} \in L_0$ . Supposons qu'il existe  $xyz = u$  tel que

$$|xy| \leq N, y \neq \varepsilon \text{ et } xy^*z \subseteq L$$

Alors, comme  $|xy| \leq N$ ,  $x$  et  $y$  sont de la forme  $a^p$  et  $a^k$ , et  $z$  est de la forme  $a^l b^{p+k+l}$ .

Donc,  $a^p(a^k)^*a^l b^{p+k+l} \subseteq L_0$ . Entre autre,  $a^{p+2k+l}b^{p+k+l} \in L_0$ . Absurde !

# Exercice

Montrer que le langage  $L_p = \{u_1 u_2 \dots u_k \mid u_1 u_2 \dots u_k = u_k \dots u_2 u_1\}$  le langage des palindromes, n'est pas rationnel.

# Exercice

Montrer que le langage  $L_p = \{u_1 u_2 \dots u_k \mid u_1 u_2 \dots u_k = u_k \dots u_2 u_1\}$  le langage des palindromes, n'est pas rationnel.

**Idée :** considérer le mot  $a^{N+1} b a^{N+1}$

# Expressions régulières en pratique

La plupart des éditeurs de texte vous permettent de rechercher des expressions régulières :

- remplacer `(\n\s*= .*)\n` par `$1.\n`

# Expressions régulières en pratique

## Syntaxe :

- `.` → n'importe quel caractère (sauf `\n`)
- `\s` → n'importe quel espace / tab / etc
- `\w` → n'importe quelle lettre
- `[abfg ]` → un des caractères parmi "a" "b" "f" "g" et " ".
- `[^ab ]` → tout sauf un des caractères parmi "a" "b" " ".
- `*`, `+`, `?`, `|` → comme ce qu'on a vous

# Expressions régulières en pratique

Remplacement :

- $\$0$  → tout
- $\$1$  → la première partie entre parenthèse
- $\$2$  → la deuxième, etc.