

I Mots qui commutent

Soient u et v deux mots. Montrer que les deux conditions suivantes sont équivalentes :

1. $uv = vu$.
2. Il existe un mot w et des entiers $k, p \in \mathbb{N}$ tels que $u = w^k$ et $v = w^p$.

Solution : Par récurrence forte sur $n = |u| + |v|$.

Cas de base : $|u| + |v| = 0$. Si $u = v = \varepsilon$, alors $u = w^1$ et $v = w^1$ pour $w = \varepsilon$, $n = p = 1$.

Cas inductif : Soit $n \in \mathbb{N}^*$. Supposons que la propriété est vraie pour des mots u, v tels que $|u| + |v| < n$.

Soient u et v tels que $uv = vu$ et $|u| + |v| < n$.

Si $|u| = |v|$ alors les $|u|$ premières lettres dans l'égalité $uv = vu$ donne $u = v$ et $u = w^1 = v$ avec $w = u$.

Supposons $|u| \leq |v|$ (l'autre cas étant symétrique). Comme $uv = vu$, u est préfixe de v : il existe un mot $v' \neq \varepsilon$ tel que $v = uv'$. On a alors $u^2v' = uv'u$. En particulier, $uv' = v'u$. Comme $|u| + |v'| < |u| + |v|$, il existe un mot w et des entiers $k, p \geq 1$ tels que $u = w^k$ et $v' = w^p$, par hypothèse de récurrence. On a alors $u = w^k$ et $v = uv' = w^{k+p}$, ce qui conclut la preuve.

II Règles sur les expressions régulières

Pour chacune des propositions suivantes sur des expressions régulières quelconques, donner une preuve ou un contre-exemple :

- | | |
|-------------------------------------|--|
| 1. $(e^*)^* \equiv e^*$ | 3. $(e_1e_2)^* \equiv e_1^*e_2^*$ |
| 2. $(e_1 e_2)^* \equiv e_1^* e_2^*$ | 4. $(e_1 e_2)^* \equiv (e_1^*e_2^*)^*$ |

Solution :

- | | |
|---|---|
| 1. Vrai. Voir cours. | 3. Faux car $abab \in (ab)^*$ mais $abab \notin a^*b^*$. |
| 2. Faux car $ab \in (a b)^*$ mais $ab \notin a^* + b^*$. | 4. Vrai. Voir cours. |

III Exemples de langages réguliers

1. Écrire une expression régulière dont le langage est l'ensemble des mots sur $\{a, b, c\}$ contenant exactement un a et un b (et un nombre quelconque de c).
2. Écrire une expression régulière dont le langage est l'ensemble des mots sur $\{a, b, c\}$ ne contenant pas de a consécutifs (aa ne doit pas apparaître).
3. Écrire une expression régulière dont le langage est l'ensemble des mots sur $\{a, b, c\}$ contenant exactement deux a et tels que tout c est précédé d'un b .
4. Si $x \in \mathbb{R}$, on note $L(x)$ l'ensemble des préfixes des chiffres de x après la virgule. Par exemple, $L(\pi) = \{\varepsilon, 1, 14, 141, 1415, \dots\}$. En sachant que $\frac{1}{6} = 0.1666\dots$ et $\frac{1}{7} = 0.142857142857\dots$, montrer que $L(\frac{1}{6})$ et $L(\frac{1}{7})$ sont réguliers.
5. Montrer plus généralement que $L(x)$ est régulier si $x \in \mathbb{Q}$ (on montrera plus tard que c'est en fait une équivalence).

Solution :

1. En distinguant le cas où a est avant b et le cas où b est avant a : $c^*ac^*bc^*|c^*bc^*ac^*$.
2. On peut donner $(a(b|c)|b|c)^*(a|\varepsilon)$ (un a doit être suivi d'un b ou d'un c).
3. Soit $e = (b|bc)^*$ (décrivant tous les mots sur $\{b, c\}$ dont chaque c est précédé d'un b). Alors $eaeae$ est une expression régulière qui convient.
4. $\varepsilon|16^*$ est une expression régulière de langage $L(\frac{1}{6})$.
 $(142857)^*(\varepsilon|1|14|142|1428|14285|142857)$ est une expression régulière de langage $L(\frac{1}{7})$.
5. Si $x \in \mathbb{Q}$, on peut écrire ses chiffres sous la forme $x = x_1.x_2ppp\dots$. Soit $Pref(m)$ l'ensemble des préfixes d'un mot m , qui est un ensemble fini si m est fini ($|Pref(m)| = |m| + 1$). Alors $L(x) = Pref(x_2)|x_2p^*Pref(p)$ (un élément de $L(x)$ est soit un préfixe de x_2 soit contient x_2 suivi d'un certain nombre

de p , suivi d'une partie de p).

IV Distance de Hamming

Si $u = u_1...u_n$ et $v = v_1...v_n$ sont deux mots de même longueur sur un alphabet Σ , leur distance de Hamming est :

$$d(u, v) = |\{i \mid u_i \neq v_i\}|$$

1. Montrer que la distance de Hamming est une distance sur Σ^* .

Solution : Soient $u = u_1...u_n, v = v_1...v_n, w = w_1...w_n$ trois mots de même taille. Si $u_i \neq w_i$ alors $u_i \neq v_i$ ou $v_i \neq w_i$ (sinon, $u_i = v_i = w_i$). D'où $d(u, v) + d(v, w) \leq d(u, w)$. $d(u, v) = d(v, u)$ et $d(u, v) = 0 \Leftrightarrow u = v$ sont facilement vérifiés.

Étant donné un langage L sur Σ , on définit son voisinage de Hamming $\mathcal{H}(L) = \{u \in \Sigma^* \mid \exists v \in L, d(u, v) = 1\}$. Pour une expression régulière e , on note $\mathcal{H}(e)$ au lieu de $\mathcal{H}(L(e))$.

2. Donner une expression régulière pour $\mathcal{H}(0^*1^*)$.

Solution : C'est l'ensemble des mots obtenus en changeant un 0 par un 1 ou inversement, c'est à dire $L(0^*10^*1^*|0^*1^*01^*)$.

3. Montrer que si L est un langage régulier alors $\mathcal{H}(L)$ est un langage régulier.

Solution : Pour simplifier, on peut raisonner par induction sur une expression régulière e de langage L .

- Si $e = \emptyset$ ou $e = \varepsilon$: $\mathcal{H}(e) = \emptyset$ est régulier.
- Si $e = a \in \Sigma$: $\mathcal{H}(e) = \Sigma \setminus \{a\}$ est régulier car fini.
- Si $e = e_1|e_2$ avec $\mathcal{H}(e_1)$ et $\mathcal{H}(e_2)$ réguliers, alors $\mathcal{H}(e) = \mathcal{H}(e_1)|\mathcal{H}(e_2)$ est régulier car union de deux langages réguliers.
- Si $e = e_1e_2$ avec $\mathcal{H}(e_1)$ et $\mathcal{H}(e_2)$ réguliers, alors $\mathcal{H}(e) = \mathcal{H}(e_1)e_2|e_1\mathcal{H}(e_2)$ est régulier car concaténation et union de langages réguliers.
- Si $e = e_1^*$ avec $\mathcal{H}(e_1)$ régulier, alors $\mathcal{H}(e) = e_1^*\mathcal{H}(e_1)e_1^*$ est régulier car étoile et concaténation de langages réguliers.

4. Écrire une fonction `h : regexp -> regexp` renvoyant une expression régulière pour le voisinage de Hamming d'un langage sur $\Sigma = \{0, 1\}$, en utilisant le type suivant :

```
type regexp =  
  | Vide | Epsilon | L of int (* L a est la lettre a (0 ou 1) *)  
  | Union of int regexp * int regexp  
  | Concat of int regexp * int regexp  
  | Etoile of int regexp
```

Solution :

```
let rec h = function  
  | Vide | Epsilon -> Vide  
  | L a -> L (1 - a)  
  | Union(e1, e2) -> Union(h e1, h e2)  
  | Concat(e1, e2) -> Union(Concat(h e1, e2), Concat(e1, h e2))  
  | Etoile e -> Concat(Etoile e, Concat(h e, Etoile e))
```

V Clôture par sur-mot (oral ENS info)

On fixe un alphabet Σ . Étant donné deux mots $w, w' \in \Sigma^*$, on dit que w' est un sur-mot de w , noté $w \preceq w'$, s'il existe une fonction strictement croissante ϕ de $\{1, \dots, |w|\}$ dans $\{1, \dots, |w'|\}$ telle que $w_i = w'_{\phi(i)}$ pour tout $1 \leq i \leq |w|$, où $|w|$ dénote la longueur de w et w_i dénote la i -ème lettre de w . Étant donné un langage L , on note \overline{L} le langage des sur-mots de mots de L , c'est-à-dire $\overline{L} := \{w' \in \Sigma^* \mid \exists w \in L, w \preceq w'\}$.

1. On pose L_0 le langage défini par l'expression régulière ab^*a , et L_1 le langage défini par l'expression régulière $(ab)^*$. Donner une expression régulière pour $\overline{L_0}$ et pour $\overline{L_1}$.
2. Montrer que, pour tout langage L , on a $\overline{\overline{L}} = \overline{L}$.
3. Existe-t-il des langages L' pour lesquels il n'existe aucun langage L tel que $\overline{L} = L'$?
4. Montrer que, pour tout langage régulier L , le langage \overline{L} est également régulier.
5. On admettra pour cette question le résultat suivant : pour toute suite $(w_n)_{n \in \mathbb{N}}$ de mots de Σ^* , il existe $i < j$ tels que $w_i \preceq w_j$.
Montrer que, pour tout langage L (non nécessairement régulier), il existe un langage fini $F \subseteq L$ tel que $\overline{F} = \overline{L}$.
6. Un langage L est clos par sur-mots si, pour tout $u \in L$ et $v \in \Sigma^*$ tel que $u \preceq v$, on a $v \in L$. Dédurre de la question précédente que tout langage clos par sur-mots est régulier.
7. On admet que les langages réguliers sont stables par passage au complémentaire. Un langage L est clos par sous-mots si, pour tout $u \in L$ et $v \in \Sigma^*$ tel que $v \preceq u$, on a $v \in L$. Montrer que tout langage clos par sous-mots est régulier.
8. Démontrer le résultat admis à la question 5.

Solution :

1. Le langage $\overline{L_0}$ est le langage des mots qui contiennent deux a , c'est-à-dire $\Sigma^*a\Sigma^*a\Sigma^*$. En effet, tout sur-mot d'un mot de ab^*a doit clairement contenir deux a . Réciproquement, tout mot contenant deux a est un sur-mot de aa qui appartient à L_0 .

Le langage $\overline{L_1}$ est Σ^* , puisque tout mot est un sur-mot de $\varepsilon \in L_1$.

2. On observe d'abord que la relation \preceq est transitive. En effet, pour tous mots $w, w', w'' \in \Sigma^*$ tels que $w \preceq w'$ et $w' \preceq w''$, en notant ϕ et ϕ' les fonctions strictement croissantes qui en témoignent, leur composition $\phi' \circ \phi$ est une fonction strictement croissante de $\{1, \dots, |w|\}$ dans $\{1, \dots, |w''|\}$, et pour tout $1 \leq i \leq |w|$ on a $w''_{\phi'(\phi(i))} = w'_{\phi(i)} = w_i$.

On montre à présent l'égalité demandée. Il est clair que $\overline{L} \subseteq \overline{\overline{L}}$, donc on montre l'inclusion inverse. Soit $u'' \in \overline{\overline{L}}$, il existe un mot $u' \in \overline{L}$ tel que $u' \preceq u''$. Par définition de \overline{L} , il existe un mot $u \in L$ tel que $u \preceq u'$. Par transitivité, on a $u \preceq u''$. Ainsi, on a bien $u'' \in \overline{L}$, ce qui conclut.

3. Pour tout langage non-vide L , le langage \overline{L} est nécessairement infini : en effet, pour $u \in L$ quelconque, on a $u\Sigma^* \subseteq \overline{L}$. Par ailleurs, on a clairement $\overline{\emptyset} = \emptyset$. Ainsi, si l'on prend L' fini non-vide, on sait qu'il n'existe aucun langage L tel que $\overline{L} = L'$.

Autre preuve possible : on considère le langage L_0 . Supposons par l'absurde qu'il existe un langage L tel que $\overline{L} = L_0$. Dans ce cas, on a $\overline{\overline{L}} = \overline{L_0}$, donc d'après la question 1, on a $\overline{L} = \overline{L_0}$. C'est absurde car L_0 et $\overline{L_0}$ sont manifestement différents. Ainsi, $L' := L_0$ convient.

4. Soit A un automate fini non-déterministe qui reconnaisse le langage régulier L . Construisons un automate A' en ajoutant à chaque état de A une boucle pour toutes les lettres de l'alphabet : formellement, on initialise $A' := A$ et pour chaque $a \in \Sigma$ et chaque état q de A , on ajoute à A' une transition de q à q étiquetée par a .

Il est clair que, pour tout mot u accepté par A et pour tout mot u' tel que $u \preceq u'$, le mot u' est accepté par A' : pour ϕ une fonction strictement croissante qui témoigne du fait que $u \preceq u'$, il suffit de suivre le chemin pour u dans A' pour les positions de u' appartenant à l'image de ϕ , et de suivre les nouvelles transitions pour les positions de u' qui n'appartiennent pas à l'image de ϕ . Réciproquement, si l'on considère un mot u' accepté par A' et un chemin qui en témoigne, on peut construire un mot u accepté par A tel que $u \preceq u'$ en considérant la restriction de ce chemin aux transitions de A .

On peut aussi démontrer cette question par induction structurelle sur les expressions régulières à l'aide des identités suivantes :

- $\bar{\emptyset} = \emptyset$
- $\bar{\varepsilon} = \Sigma^*$
- $\bar{a} = \Sigma^* a \Sigma^*$ pour tout $a \in \Sigma$ $\overline{L_1 L_2} = \overline{L_1} \overline{L_2}$
- $\overline{L_1 \cup L_2} = \overline{L_1} \cup \overline{L_2}$
- $\overline{L^*} = \Sigma^*$

La dernière égalité est due au fait que L^* contient toujours le mot vide; en revanche il n'est pas vrai que $\overline{L^*} = \overline{L}^*$, prendre par exemple $L = a$.

5. Soit L un langage quelconque. Si L est vide, on peut prendre $F = \emptyset$ et conclure. Sinon, posons $(w_n)_{n \in \mathbb{N}}$ une suite infinie énumérant les mots du langage L (éventuellement avec des doublons). Une position $i \in \mathbb{N}$ est dite innovante s'il n'existe aucun $j < i$ tel que $w_j \preccurlyeq w_i$. On choisit pour F le sous-ensemble de L formé des mots aux positions innovantes, c'est-à-dire $\{w_i \mid i \text{ est innovante}\}$.

On observe à présent qu'il y a un nombre fini de positions innovantes. En effet, dans le cas contraire, la suite extraite obtenue à partir de $(w_n)_{n \in \mathbb{N}}$ en conservant les lettres aux positions innovantes serait un contre-exemple à la question 4. Ainsi, F est-il bien fini.

Montrons à présent que $\overline{F} = \overline{L}$. En effet, comme $F \subseteq L$, on a $\overline{F} \subseteq \overline{L}$ par monotonie de la clôture par sur-mots. Pour la réciproque, il suffit de montrer que $L \subseteq \overline{F}$, car cela implique (à nouveau par monotonie de la clôture par sur-mots) que $\overline{L} \subseteq \overline{\overline{F}}$, ce qui implique par la question 1 que $\overline{L} \subseteq \overline{F}$. Montrons par induction sur $i \in \mathbb{N}$ que $w_j \in \overline{F}$ pour tout $j < i$. Le cas de base est tautologique. Pour le cas de récurrence, choisissons $i \in \mathbb{N}$. Soit i est innovante, soit i n'est pas innovante. Dans le premier cas, on a $w_i \in F$ donc $w_i \in \overline{F}$. Dans le second cas, il existe $j < i$ tel que $w_j \preccurlyeq w_i$, et par hypothèse de récurrence on a $w_j \in \overline{F}$, ainsi on a $w_i \in \overline{F}$. Ainsi, dans les deux cas on a $w_i \in \overline{F}$. On a donc établi notre résultat par récurrence, et on a donc bien l'inclusion réciproque $L \subseteq \overline{F}$.

6. Soit L un langage clos par sur-mots. On sait par la question 5 qu'il existe un langage fini $F \subseteq L$ tel que $\overline{F} = \overline{L}$. Or on a $\overline{L} = L$. En effet, il est clair que $L \subseteq \overline{L}$, et réciproquement, pour tout $v \in \overline{L}$, il existe par définition de \overline{L} un mot $u \in L$ tel que $u \preccurlyeq v$, et ainsi $v \in L$ car L est clos par sur-mots. On sait donc que $L = \overline{F}$, et on sait que F est régulier (car fini), donc \overline{F} est régulier par la question 3, ainsi L est-il régulier.
7. Soit L un langage clos par sous-mots, et soit $L' := \Sigma^* \setminus L$ son complémentaire. Montrons que L' est clos par sur-mots. En effet, soit $u \in L'$ et $v \in \Sigma^*$ tels que $u \preccurlyeq v$. Procédons par l'absurde et supposons que $v \notin L'$. On a alors $v \in L$. Comme $u \preccurlyeq v$ et que L est clos par sous-mots, on sait que $u \in L$, et ainsi $u \notin L'$, contredisant notre hypothèse. Ainsi, $v \in L'$, ce qui établit que L' est clos par sur-mots. On sait donc que L' est régulier. Comme les langages réguliers sont clos par complémentation, le complémentaire L de L' est lui aussi régulier.
8. Il s'agit du lemme de Higman dans le cas particulier des alphabets finis.

Procédons par l'absurde et supposons qu'il existe une mauvaise suite, c'est-à-dire une suite $(w_n)_{n \in \mathbb{N}}$ telle qu'il n'existe pas de $i < j$ telle que $w_i \preccurlyeq w_j$. Construisons une nouvelle suite $(w'_n)_{n \in \mathbb{N}}$ de la façon suivante : le mot w'_0 est un mot de longueur minimale telle qu'il existe une mauvaise suite commençant par w'_0 (un tel w'_0 existe par notre hypothèse), le mot w'_1 est un mot de longueur minimale telle qu'il existe une mauvaise suite commençant par w'_0, w'_1 (un tel w'_1 existe par définition de w'_0), et ainsi de suite. La suite $(w'_n)_{n \in \mathbb{N}}$ ainsi définie est clairement mauvaise : pour tout $i < j$, la définition de w'_j assure qu'on ne peut avoir $w'_i \preccurlyeq w'_j$. Par ailleurs, la définition de $(w'_n)_{n \in \mathbb{N}}$ assure qu'elle est minimale, c'est-à-dire que pour toute mauvaise suite $(w''_n)_{n \in \mathbb{N}}$, si on pose $i \in \mathbb{N}$ le premier indice tel que $w'_i \neq w''_i$, on a nécessairement $|w'_i| \leq |w''_i|$.

On va aboutir à notre contradiction en construisant à partir de $(w'_n)_{n \in \mathbb{N}}$ une nouvelle mauvaise suite qui contredise sa minimalité. Soit $a \in \Sigma$ une lettre quelconque telle qu'il existe un nombre infini de mots de $(w'_n)_{n \in \mathbb{N}}$ ayant a comme première lettre : comme Σ est fini, un tel a existe nécessairement. Soit $p \in \mathbb{N}$ le plus petit entier tel que w'_p commence par a . On construit la suite $(w''_n)_{n \in \mathbb{N}}$ comme la concaténation de w'_0, \dots, w'_{p-1} et de la suite extraite de $(w'_n)_{n \in \mathbb{N}}$ des mots commençant par a à qui on a retiré leur première lettre.

La suite $(w''_n)_{n \in \mathbb{N}}$ est une mauvaise suite. En effet, soit $p < q$. Si $q < i$, alors $w''_p = w'_p$ et $w''_q = w'_q$, donc $w''_p \not\preccurlyeq w''_q$ car $(w'_n)_{n \in \mathbb{N}}$ est mauvaise. Si $i \leq p$, alors $aw''_p = w'_p$ et $aw''_q = w'_q$, donc on conclut encore car $(w'_n)_{n \in \mathbb{N}}$ est mauvaise. Si $p < i \leq q$, alors $w''_p = w'_p$ et $aw''_q = w'_q$, et on conclut de même.

Par ailleurs, la suite $(w''_n)_{n \in \mathbb{N}}$ contredit la minimalité de $(w'_n)_{n \in \mathbb{N}}$. En effet, le premier indice où ces deux suites diffèrent est p , et on a $|w''_p| = |w'_p| - 1$, ce qui contredit bien la minimalité. C'est absurde, et ainsi notre hypothèse initiale affirmant l'existence d'une mauvaise suite est-elle fausse.