

SÉMINAIRE DUBREIL. ALGÈBRE ET THÉORIE DES NOMBRES

M. P. SCHÜTZENBERGER

Une théorie algébrique du codage

Séminaire Dubreil. Algèbre et théorie des nombres, tome 9 (1955-1956), exp. n° 15,
p. 1-24

http://www.numdam.org/item?id=SD_1955-1956__9__A10_0

© Séminaire Dubreil. Algèbre et théorie des nombres
(Secrétariat mathématique, Paris), 1955-1956, tous droits réservés.

L'accès aux archives de la collection « Séminaire Dubreil. Algèbre et théorie des nombres » implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques
<http://www.numdam.org/>

UNE THÉORIE ALGÈBRIQUE DU CODAGE,

par M.P. SCHÜTZENBERGER.

Introduction.

Soient donnés deux ensembles discrets (en général finis dans la pratique) :

Λ_0 : l'ensemble des messages élémentaires

A_0 : l'ensemble des lettres ou "alphabet".

En théorie des communications, le problème du codage consiste à représenter les messages, c'est-à-dire les suites de messages élémentaires par des suites de lettres selon une règle fixée à l'avance (un "code" ou "dictionnaire") de telle sorte que la retraduction soit possible et ceci en fonction d'exigences techniques de nature diverse.

Par exemple, Λ_0 étant l'ensemble des signes typographiques habituels et A_0 un "alphabet" consistant en les trois "lettres" "point" "trait" et "intervalle".

Le code télégraphique morse est le système de convention dans lequel "a" est représenté par "point trait intervalle", "b" par "trait point point point intervalle" ... etc. On appellera respectivement "codage" et "décodage" les deux opérations inverses qui font passer des suites de messages aux suites de lettres et réciproquement.

Plus généralement on pourrait voir dans tout langage, naturel ou artificiel, un code (plus souvent d'ailleurs : une série de codes superposés : phonémiques, grammaticaux, sémantiques, etc.) traduisant en sons, signes ou gestes des idées, des sentiments, etc. Une définition encore plus générale a été étudiée par J. Riguet [1], et nous désignerons ici, pour abrégé, sous le nom de "codes", certaines structures qui idéalisent les restrictions plus ou moins explicitement ou rigoureusement admises par la quasi-totalité des systèmes utilisés dans la théorie des communications stricto sensu :

- 1°) Universalité : Toute suite de messages élémentaires correspond à une suite de lettres au moins.

- 2°) Catégoricité : Cette suite de lettre est unique pour une suite de messages élémentaires donnée.
- 3°) Univocité du décodage : Réciproquement, si une suite de lettres correspond à une suite de messages élémentaires, cette dernière est unique.
- 4°) Isomorphie : Si les messages " λ " et " μ " sont codés respectivement par " ℓ " et " m ", le message " $\lambda\mu$ " est codé par " ℓm ".

Le code morse satisfait rigoureusement à ces quatre conditions et aussi, dans une certaine mesure, ce code compliqué qu'est la transcription orthographique française du langage parlé⁽¹⁾.

Formellement il est clair que les conditions énoncées conduisent à la :

Définition :

Λ et A étant respectivement les demi-groupes libres engendrés par Λ_0 et A_0 , une application \mathcal{C} de Λ_0 sur une partie P_0 de A sera un code, si et seulement si l'extension de \mathcal{C} à Λ détermine un isomorphisme de Λ sur le sous-demi-groupe P de A engendré par les suites de lettres constituant P_0 .

On appellera "mot" les suites de lettres constituant P_0 et, pour abrégé, "message" aussi bien les éléments de P que ceux de Λ . Un code sera défini par la donnée de Λ_0 , A_0 et \mathcal{C} , ou plus fréquemment par celle de A et de P , puisque, A et P étant des demi-groupes, leur donnée implique celle de leurs générateurs.

(1) Les contre-exemples suivants illustrent la signification des conditions précédentes :

- ad 1°) : le phonème " t^* " (le click dental du Bantou, utilisé en France pour stimuler les chevaux) est intranscriptible dans l'orthographe française.
- ad 2°) : le phonème "K" est codé K (Képi) ou ch (choléra) ou c (carotte) ou qu (quai).
- ad 3°) : la graphie "ent" code plusieurs phonèmes distincts (vent, vient, viennent).
- ad 4°) : la graphie "oi" qui se lit "ua" et devrait s'écrire "oua".

La définition précédente fait apparaître la théorie de codage comme une application de la théorie des demi-groupes. Il est particulièrement remarquable que les concepts fondamentaux de cette dernière, introduits par P. Dubreil [3], [4], [5] en 1941, et étudiés depuis par lui-même et son école du point de vue abstrait, aient des interprétations immédiates et importantes sur le plan de la réalisation concrète des machines codeuses ou transcodeuses.

Sous un autre angle, la théorie des événements récurrents de W. Feller [6], dans laquelle on étudie des processus stochastiques sur les suites de lettres à partir d'autres processus définis sur les suites de messages élémentaires, et réciproquement se rattache étroitement à la théorie du codage, comme l'a montré Mandelbrot qui a utilisé cette analogie dès 1951 [7], [8] : les processus récurrents sont des codes unitaires au sens que nous donnerons plus loin à ce terme. En retour, l'extension par Feller lui-même de sa théorie à des alphabets topologiques spéciaux, ouvre la voie à des généralisations intéressantes qui feront l'objet d'un exposé ultérieur.

Enfin, la théorie développée ici est une théorie "sans bruit" dans laquelle on ne se préoccupe pas des problèmes que pose au décodeur une altération d'origine aléatoire des messages à décoder. Le contraste avec les travaux actuels sur cette question (Cf: [10] ou [11] par exemple) est moins grand qu'il ne paraît. En particulier, la notion de "code absorbant" y joue un rôle considérable, et ne sera ici l'objet que de remarques incidentes.

Exemple.

L'exposé sera simplifié si un exemple est traité qui montre la nature du problème :

Soit le code ⁽²⁾ suivant dans l'alphabet $A_0 = \{ a, b, c \}$:

$$C_\alpha \rightarrow a ; \quad C_\beta = bb ; \quad C_\gamma = c ; \quad C_\delta = ab ; \quad C_\varepsilon = bcb .$$

Il n'y a aucune difficulté à réaliser le codage d'une suite de messages élémentaires. Par exemple : $C(\alpha \varepsilon \delta \varepsilon) = abcbabbcb = m$.

Le décodage par contre n'est pas aussi direct : en procédant de gauche à droite, en effet, le "a" initial est, a priori, aussi bien " α " que la première lettre du mot ε . Cependant cette dernière hypothèse est prouvée fausse, car "c" (troisième lettre du message) devrait être lu " γ " et la séquence restante (babbcb) ne pourrait pas être décodée puisqu'aucun mot ne commence par "ba". Donc selon des notations évidentes :

(2) On prouvera plus loin qu'il s'agit bien d'un code.

$$m = \alpha . bcb.abbcb = \alpha \xi abbcb = \alpha \xi m'$$

La même ambiguïté se présente à nouveau : si la lettre initiale "a" du message restant m' était α on aurait nécessairement : $m' = \alpha bcb = \alpha \beta cb = \alpha \beta \gamma b$ et b resterait sans pouvoir être décodé. Donc : $m' = \delta bcb = \delta \xi$ et l'on retrouve bien $G_m^{-1} = \alpha \xi \delta \xi$. Observons au passage qu'il se peut que l'ambiguïté ne soit susceptible d'être levée qu'avant un délai très long :

Le message $q = a(bbc)^n$ (a suivi de n-fois bbc) correspond à $\alpha (\beta \gamma)^n$, et le message $qb = a(bbc)^n b$ correspond à $\delta \xi^n$.

Il est donc nécessaire :

- 1°) De systématiser les méthodes de décodage (2e section).
- 2°) De rendre possible leur mécanisation en les simplifiant au maximum par l'emploi de structures moins lourdes que les demi-groupes libres A et P : (3e section).
- 3°) D'étudier des codes particulièrement maniables (les codes unitaires nets) et de montrer qu'ils forment une "classe admissible" en ce sens qu'ils sont aussi efficaces que n'importe quel autre code du point de vue de la longueur moyenne des mots. (3e et 4e section).

1.- Méthodes de décodage.

Soient A un demi-groupe libre, P un sous-demi-groupe de celui-ci. Le premier problème sera de caractériser algébriquement les P susceptibles de correspondre à un code.

On posera : $P_0 = P - ((P - \emptyset)^2 \cup \emptyset) =$ l'ensemble des $p \in P - \emptyset$ tels que $p'p'' \in P$ et $p', p'' \in P$ impliquent $p' = \emptyset$ ou $p'' = \emptyset$ (3). Comme A est un demi-groupe libre chaque élément possède une "longueur" (le nombre de générateurs figurant dans son expression) et par récurrence, il est clair que tout $p \in P$ est d'au moins une manière décomposable en un produit $p = p_1 p_2 \dots p_m$ où tous les p_i appartiennent à P_0 .

Proposition 1.1.- Une condition nécessaire et suffisante pour que $P \subset A$ corresponde à un code est que la décomposition en produits d'éléments $p_i \in P_0$ soit unique pour tout $p \in P$.

(3) Ici, et par la suite, on désignera par \emptyset la suite vide comme distincte de \emptyset : l'ensemble vide. On conviendra en outre que la "suite vide" est un élément de tout demi-groupe considéré dont elle est évidemment un élément neutre.

Attachons en effet à tout p_i un $\lambda_i = \mathcal{C}_{p_i}^{-1}$ par une correspondance biunivoque. L'extension de \mathcal{C} au demi-groupe libre Λ engendré par les λ_i est un homomorphisme de Λ sur P puisque si $\lambda = \lambda_1 \lambda_2 \dots \lambda_m$,
 $\mathcal{C}\lambda = \mathcal{C}\lambda_1 \mathcal{C}\lambda_2 \dots \mathcal{C}\lambda_m = p_1 p_2 \dots p_m$. Cet homomorphisme est un isomorphisme sur si et seulement si $\lambda \neq \lambda'$ entraîne $\mathcal{C}\lambda \neq \mathcal{C}\lambda'$ c'est-à-dire si et seulement si $p_1 p_2 \dots p_m \neq p'_1 p'_2 \dots p'_m$ quels que soient $p_i, p'_i \in P_0$.

Considérons, pour P quelconque, un élément $a = a_1 a_2 \dots a_n$ de \mathbb{A} où les a_i sont des lettres ($a_i \in A_0$).

Définition. On appellera indice critique de a tout indice i tel que $a_1 a_2 \dots a_i = {}^*a_i \in P$ et $a_{i+1} a_{i+2} \dots a_n = a_i^* \in P$. Evidemment l'ensemble J_a des indices critiques de a n'est non vide que si $a = {}^*a_i a_i^* \in P$.

Proposition 1.2.— Si J est un indice critique de la sous-séquence de a $b = b_{i,i'} = a_{i+1} a_{i+2} \dots a_{i'}$, où $i, i' \in J_a$ il est aussi un indice critique de a .

En effet si ${}^*b_J = a_{i+1} a_{i+2} \dots a_J \in P$ et $b_J^* = a_{J+1} \dots a_{i'} \in P$ on a aussi ${}^*a_i {}^*b_J = {}^*a_J \in P$ et $b_J^* a_{i'}^* = a_J^* \in P$.

Proposition 1.3.— Une condition nécessaire et suffisante pour que $P \subset A$ corresponde à un code est que, pour tout $a \in P$ et toute sous-séquence de a , $b = b_{i,i'}$, $i, i' \in J_a$ entraîne $b \in P$.

La condition est suffisante car si $J_a = \{1 = i_1 < i_2 < \dots < i_m = n\}$ d'une part les sous-séquences $b_{i_k i_{k+1}}$ sont indécomposables d'après 1:2 et sont donc des mots,

d'autre part il ne peut exister aucune autre décomposition de a puisque celle-ci impliquerait l'existence d'un indice critique $i' \notin J_a$.

La condition est nécessaire : supposons que $b_{i,i'} \notin P$. Il existe dans J_a J et J' (éventuellement : $J = 1$ et $J' = n$) tels que :

$$1^\circ) \quad J < i < i' < J'$$

$$2^\circ) \quad b_{Ji}, b_{Ji'}, b_{iJ'}, b_{JJ'} \in P$$

$$3^\circ) \quad |J' - J|$$

soit minimum parmi les couples J et J' satisfaisant à 1°) et 2°) .

$c = b_{JJ'}$ possède deux décompositions en mots au moins : l'une avec :

$$p_1 p_2 \dots p_m' = b_{Ji'} \quad \text{et} \quad p_{m'+1} \dots p_m = b_{i'J'}$$

l'autre avec

$$p_1' p_2' \dots p_m' = b_{Ji} \quad \text{et} \quad p_{m'+1}' \dots p_m' = b_{iJ'} \quad (p_i \in P_0)$$

Ces deux décompositions sont distinctes, car si par exemple on avait $p_1 = p'_1 = a_1 a_2 \dots a_{J''}$ le couple (J'', J') satisferait à 1°) et 2°) et on aurait $|J' - J''| < |J' - J|$ en contradiction avec 3°). Une conséquence de la proposition précédente est le

Théorème 1.4. - Une condition nécessaire et suffisante pour que le sous-demi-groupe P du demi-groupe libre A corresponde à un groupe est que :

$$(1.4) \quad P^{(-1)} P \cap P P^{(-1)} \subset P \quad (4)$$

En effet la condition indiquée signifie simplement que : $p, p', pq, qp' \in P$ implique $q \in P$ et il suffit d'appliquer la proposition en prenant $a = pqp' \in P$.

Si $P \subset A$, où A n'est pas nécessairement un demi-groupe libre, satisfait 1.4, on dira que P est "libérable" dans A . Un autre corollaire intéressant est obtenu en remarquant que $P \subset A$ correspond à un code si et seulement si P est isomorphe à un demi-groupe libre. Soit \mathcal{L}_A l'ensemble des sous-demi-groupe de A qui satisfont à cette propriété et soit δ le "demi-groupe libre vide" ($P_0 = \emptyset$) que l'on supposera appartenir à \mathcal{L}_A :

Proposition 1.5. - \mathcal{L}_A est un treillis.

Comme $A \in \mathcal{L}_A$, il suffit de montrer que $P \in \mathcal{L}_A$ et $P' \in \mathcal{L}_A$ entraîne $P \cap P' \in \mathcal{L}_A$ où $P \cap P' = P'' \in \mathcal{L}_A$. Or si $p_1 p_2, p_1 q, qp_2 \in P''$, $q \in P$ et $q \in P'$ puisque P et P' satisfont à la relation 1.4, donc $q \in P''$ et P'' satisfait aussi 1.4.

Définition. Un code sera dit unitaire à gauche si $P^{(-1)} P \subset P$ (unitaire à droite si $PP^{(-1)} \subset P$)

(4) Les conventions d'écritures suivantes seront systématiquement employées : (à droite ou à gauche).

$$X^{[-1]} Y = Y \cdot X = \{ \hat{z} : \bigvee_X x \quad xz \in Y \}$$

$$X^{(-1)} Y = A - (A - Y) \cdot X = \{ \hat{z} : \bigwedge_X x \quad xz \in Y \}$$

si X est un élément, les symboles $x^{[-1]}$ et $x^{(-1)}$ ont le même sens. On observera que les "parenthèses" satisfont à une série d'identité parallèle à celle qui est bien connue pour les "crochets". Notamment :

$$X^{(-1)} (Y^{(-1)} Z) = (YX)^{(-1)} Z \quad ; \quad X^{(-1)} (ZY^{(-1)}) = (X^{(-1)} Z) Y^{(-1)} ; \\ X^{(-1)} (XY) \supset Y \text{ etc. Toutefois :}$$

$$(X \cup Y)^{(-1)} Z = X^{(-1)} Z \cup Y^{(-1)} Z \quad \text{et} \quad (X \cup Y)^{[-1]} Z = X^{[-1]} Z \cap Y^{[-1]} Z, \text{ etc.}$$

On écrira $X^{(-n)} Z$ pour $X^{(-1)} (X^{(-1)} (X^{(-1)} \dots Z)) = (X^n)^{(-1)} Z$.

On a évidemment :

Proposition 1.6.- Tout sous-demi-groupe unitaire P d'un demi-groupe libre correspond à un code.

La définition précédente est exactement celle de P. Dubreil [3]. Elle conduit à un décodage très simple : soit en effet $a = a_1 \dots a_n \in P$ une suite de lettres. Procédant de gauche à droite, soit $i \neq 1$ le plus petit indice tel que ${}^*a_i \in P$. Par hypothèse, puisque P est unitaire, $a_i^* \in P$ et $i \in J_a$. Donc $b_{1i} \in P_0$ et il suffit de recommencer sur la suite $a' = a_{i+1} \dots a_n \in P$ pour aboutir au décodage de a .

Exemple : Soit le code "opposé" de celui donné en exemple dans l'introduction :

$$\mathcal{C}_\alpha = a ; \mathcal{C}_\beta = bb ; \mathcal{C}_\gamma = c ; \mathcal{C}_\delta = ba ; \mathcal{C}_\varepsilon = bcb$$

On vérifie (Cf. plus bas) que P est unitaire et que le décodage de $m = bcbabcbcb$ s'effectue directement sans essais et erreurs :

$$m = \varepsilon babcba = \varepsilon \delta bcbcb = \varepsilon \delta \varepsilon a = \varepsilon \varepsilon \alpha.$$

Le problème consistant à déterminer pratiquement si un ensemble P_0 de mots correspond ou non à un code avait été résolu par Sardinas et Patterson [12] par des considérations purement combinatoires sans faire appel à la notion de demi-groupe. Leurs résultats peuvent être améliorés de la façon suivante :

supposons que $\emptyset \neq P_0$ et que $P_0 \cap P_0^2 = \emptyset$

Proposition 1.7.- Si l'on pose $P_1 = P_0^{(-1)} P_0$ et par récurrence :

$$P_{n+1} = P_n^{(-1)} P_0 \cup P_0^{(-1)} P_n \text{ alors : } P_n = \bigcup_{n'+n''=n} P_0^{-(n')} P_0^{n''}.$$

Par définition $P_1 =$ l'ensemble des $q \in A$ tels qu'il existe $p, p' \in P_0$ avec $pq = p'$. Supposons établi que pour $n \leq n_0$ on ait montré que :

$P_n =$ l'ensemble des q tels qu'il existe $p_1, p_2 \dots p_{n'}$, $p_1' p_2' \dots p_{n''}' \in P_0$ ($n' + n'' = n$) avec

$$p_1 p_2 \dots p_{n'} q = p_1' p_2' \dots p_{n''}' \quad (1)$$

Tout élément $q' \in P_{n+1}$ est défini par :

$$\text{soit } pq' = q \quad (2)$$

$$\text{soit } qq' = p \quad (3)$$

avec $p \in P_0$, $q \in P_n$

Dans le premier cas, multiplions (2) à gauche par $p_1 p_2 \dots p_n$, il vient :

$$p_1 p_2 \dots p_n p q' = p_1 p_2 \dots p_n q \quad \text{soit : } q' = P^{-(n+1)} P^n .$$

Dans le second cas, multiplions (1) à droite par q' , il vient :

$$p_1 p_2 \dots p_n q q' = p'_1 p'_2 \dots p'_n q' \quad \text{soit d'après (3) : } q' = P^{(-n)} P^{n+1} .$$

Proposition 1.8. - Une condition nécessaire et suffisante pour que P_0 soit l'ensemble des mots d'un code est que, pour tout $n \geq 1$, $P_n \cap P_0 = \emptyset$.

En effet, on aura ou non affaire à un code selon que pour tout $a \in A$

$$a = p_1 p_2 \dots p_m = p'_1 p'_2 \dots p'_m q \quad (p_i, p'_i \in P_0 ; p_i \neq p'_i)$$

entraînera ou non $q \notin P$.

L'intérêt de cette proposition est que, si le code est borné c'est-à-dire si la longueur de tous ses mots est bornée, les $q \in P_n$ sont des diviseurs à gauche ou à droite des séquences composant P_0 , donc sont bornés eux-mêmes et en nombre fini s'il en est de même de la puissance de l'alphabet. Dans ces conditions il existe un $n < \infty$ tel que

$$\text{soit } P_n = \emptyset \quad \text{et par conséquent } P_{n'} = \emptyset \quad \text{pour tout } n' \geq n$$

soit $P_n = P_{n+r}$ et par conséquent $P_{n+J+kr} = P_{n+J}$ pour tout $J \leq r$ et tout k .

On a donc · Dans un code borné

Proposition 1.9. - Une condition nécessaire et suffisante pour qu'il existe une valeur fixe $L < \infty$ telle que la connaissance des $m + L$ premières lettres (à gauche) du message permette quel que soit m de décoder les m premières lettres sans ambiguïté est qu'il existe un $n < \infty$ tel que $P_n = \emptyset$.

En particulier une condition nécessaire et suffisante pour que le code soit unitaire à gauche est que $P_1 = \emptyset$.

En effet si un tel n n'existe pas, on peut construire des $a \in A$ de la forme 1.9 pour des valeurs aussi grandes que l'on veut de $m + m'$. Inversement si $P_n = \emptyset$ et si $a = p_1 \dots p_m = p'_1 p'_2 \dots p'_m q$ on a $p_i = p'_i$ pour tout $i \leq m_0$ avec $m - m_0 + m' - m_0 = n$.

Les notions suivantes sont utiles pour caractériser certains types de codes :

Définition. Un code sera dit

net (à droite) si $PA^{(-1)} = A$

absorbant (à droite) si $PP^{(-1)} = A$

La notion de code net est encore exactement celle de P. Dubreil [3]: un code est net à droite si quel que soit la séquence de lettres $a = a_1 a_2 \dots a_n \notin P$ il existe au moins une séquence $b = a_{n+1} a_{n+2} \dots a_m$ tel que $ab \in P$.

Le code orthographique du français n'est pas net, car il est impossible d'ajouter des lettres à droite de la séquence "Khtg" par exemple qui la complète en un mot. Par contre le code formé par l'ensemble P des phrases sémantiquement correctes est net, car à toute suite de mots ou de signes on peut au besoin rajouter la clause. "Les x-dernières syllabes que j'ai énoncées étaient un exemple de sentence sémantiquement absurde."

La notion de code absorbant est encore plus forte que la notion de code net : elle signifie que quel que soit $a \notin P$ il existe $q \in P$ (et non pas seulement dans A) tel que $aq \in P$. On démontre :

Proposition 1.10. - Si A est fini, quel que soit $P \subset A$, $A^{[-1]}_P \neq \emptyset$ est équivalent à $PP^{(-1)} = A$.

En effet $A^{[-1]}_P \ni r$ signifie que $ar \in P$ pour tout $a \in A$.

Réciproquement soit $p_1 \in P$ tel $a_1 p_1 \in P$ pour au moins un $a_1 \in A - P = A_1$ et par récurrence :

$p_i \in P$ tel que $a_i p_i \in P$ pour au moins un $a_i \in A_i = A_{i-1} p_{i-1}^{-1} - P$.

Les ensembles $A'_1 = P p_1^{-1}$; $A'_2 = P(p_1 p_2)^{-1}$ soient strictement croissants ; donc pour un certain $i = J$: $A_J = \emptyset$ et $p_1 p_2 \dots p_J \in A^{[-1]}_P$. Il est clair que si un code est unitaire net ou absorbant à droite, il ne s'en suit pas qu'il le soit aussi à gauche. En particulier :

Proposition 1.11. - Une condition nécessaire (mais non suffisante) pour qu'un code soit absorbant à droite est qu'il soit unitaire à gauche et net à droite et qu'il ne soit pas unitaire à droite.

En effet si $PP^{(-1)} = A$, P n'est libérable ($P^{(-1)} P \cap PP^{(-1)} \subset P$) que s'il est unitaire à gauche. D'autre part s'il n'était pas net, il existerait w tel que $wx \notin P$ pour tout x .

Remarque. Les énoncés précédents semblent avoir une partie pratique assez faible puisque, par exemple la recherche des indices critiques d'une séquence s'effectue pratiquement en vérifiant que $*a_i$ et a_i^* sont décodables. Pour aller plus loin il nous faudrait trouver des demi-groupes \bar{A} et \bar{P} de préférence finis et une application φ tels que pour tout $a \in A$, $\varphi a \in \bar{P}$

entraîne $a \in P$. C'est ce que nous ferons dans la section suivante. Il se trouve cependant que ces notions dépassent largement le cadre des problèmes de codage et il a paru aussi simple de les traiter en toute généralité pour un A quelconque et un complexe $K \subset A$ qui n'en est pas nécessairement un sous-demi-groupe.

2.- Equivalence syntaxiques.

Soient A un demi-groupe contenant un élément neutre, K une partie quelconque de A .

Définition. On dira que a est syntactiquement plus fort que b dans A , par rapport à K ($a \succcurlyeq b (A, K)$) si pour tout $x, y \in A$:

$$(II) \quad xby \in K \text{ entraîne } xay \in K$$

Si $a \succcurlyeq b (A, K)$ et $b \succcurlyeq c (A, K)$, a et c seront "syntactiquement équivalents" ($a \equiv c (A, K)$).

Proposition 2.1.- $\succcurlyeq (A, K)$ est une relation régulière, réflexive, transitive (une relation de préordre (2) compatible avec la structure de demi-groupe). Manifestement $a \succcurlyeq a (A, K)$ et $a \succcurlyeq b (A, K)$ et $b \succcurlyeq c (A, K)$ entraînent $a \succcurlyeq c (A, K)$.

D'autre part, si (II) est vraie pour tout $x, y \in A$, (II) est encore vraie pour $x \in Av$ et $y \in vA$. Donc $a \succcurlyeq b (A, K)$ entraîne $uav \succcurlyeq ubv (A, K)$ et en particulier si $a \succcurlyeq a' (A, K)$ et $b \succcurlyeq b' (A, K)$ on a successivement:

$$ab \succcurlyeq a'b (A, K) \text{ et } a'b \succcurlyeq a'b' (A, K) \text{ d'où } aa' \succcurlyeq bb' (A, K).$$

Proposition 2.2.- $\succcurlyeq (A, K)$ est la plus forte des relations régulières de préordre pour laquelle K soit supérieurement saturé ⁽⁵⁾. K est supérieurement saturé pour $\succcurlyeq (A, K)$ car $b \in K$ et $a \succcurlyeq b (A, K)$ entraîne en particulier $a \in K$ en faisant $x = y = \emptyset$ dans (II). Réciproquement, si ρ est régulière et K supérieurement saturé pour ρ , $a \rho b$ implique $(xay) \rho (xby)$ et en particulier $xby \in K$ entraîne $xay \in K$ donc $a \rho b$ implique $a \succcurlyeq b (A, K)$.

Proposition 2.3.- Si A est un groupe fini $\succcurlyeq (A, K)$ se réduit à une relation d'équivalence qui est précisément l'équivalence normale associée au plus grand sous-groupe normal G de A qui soit contenu dans l'un des complexes de K .

(5) ρ étant une relation de préordre (\succcurlyeq réflexive et transitive) on dira que K est supérieurement saturé pour ρ si $a \rho b$ et $b \in K$ entraînent $a \in K$. ρ est régulière si $a \rho b$ entraîne $(xay) \rho (xby)$ pour tout x, y . (Cf [3]).

En effet : si A est un groupe $a \geq b (A, K)$ peut s'écrire :
 $K y^{-1} a^{-1} b y \subset K$ c'est-à-dire encore $a^{-1} b \in G$ où G est l'ensemble des c
 tels que $K y^{-1} c y \subset K$ pour tout y . Mais, d'une part $K d \subset K$ est équivalent
 à $K d = K$, d'autre part $c, c' \in G$ entraîne $cc' \in G$. Donc G est le plus
 grand sous-groupe normal de A tel que $KG = K$ et l'on a : $K = \bigcup_{x \in K} xG$
 c'est-à-dire $G \subset k^{-1}K$ pour tout $k \in K$.

Remarque 1.- Si l'on voulait interpréter l'équivalence syntaxique dans le code
 des "phrases françaises grammaticalement correctes" on trouverait par exemple
 que "postulat" \neq "axiome" \equiv "hippopotame" car on peut aussi bien dire
 "l'axiome d'Euclide est à la base de la géométrie élémentaire" que "l'hippopo-
 tame d'Euclide ...etc" alors que la phrase "l'postulat d'Euclide ...etc" est
 incorrecte.

Remarque 2.- Dans ce même cadre linguistique il serait possible et utile de
 généraliser la définition précédente

Définition.- Le n -tuple $a = (a_1, \dots, a_n)$ d'éléments de A est "syntaxique-
 ment plus fort que le n -tuple $b = (b_1, b_2, \dots, b_n)$ dans A par rapport à K
 et au $n+1$ -tuple $(g_0, g_1, \dots, g_n) = g$ " si pour tout $2n$ -tuple
 $(x_1, y_1, x_2, y_2, \dots, x_n, y_n) = (x, y)$ $g_0 x_1 b_1 y_1 g_1 x_2 b_2 \dots g_{n-1} x_n b_n y_n g_n =$
 $= g(x, y) (b) \in K$ entraîne $g(x, y) (a) \in K$. ($a \geq b (A, K; g)$)

Les propriétés 2.1 et 2.2 subsistent et naturellement $a_i \equiv b_i (A, K)$
 pour tout i entraîne $a \geq b (A, K; g)$ la réciproque n'étant pas forcément
 vraie.

Nous ne ferons pas usage ici de cette notion générale.

Définition. Soit φ_K l'homomorphisme attaché à $\equiv (A, K)$. On posera
 $\bar{A} = \varphi_K A$; $\bar{K} = \varphi_K K$; si $K = P \supset P^2$ on dira que A et P sont les demi-
groupes syntaxiques fondamentaux (GSF) de $(A \supset P)$.

Si $\varphi_K = 1$ (φ_K réduit à l'application identique) on dira que
 $K (= \bar{K})$ est syntaxiquement simple dans $A (= \bar{A})$.

Proposition 2.4.- $a \geq b (A, K)$ est équivalent à $\varphi_K a \geq \varphi_K b (\bar{A}, \bar{K})$.
 Donc \bar{K} est syntaxiquement simple dans \bar{A} .

D'une part $xy \in K$ entraîne $\varphi_K x \varphi_K b \varphi_K y \in \varphi_K K$.

D'autre part : $\varphi_K x \varphi_K b \varphi_K y = \varphi_K (xby) \in \varphi_K K$ entraîne $xby \in K$ puis-
 que K est saturé. Donc $\varphi_K a \equiv \varphi_K b (\varphi_K A; \varphi_K K)$ est équivalent à
 $a \equiv b (A, K)$ donc encore à $\varphi_K a = \varphi_K b$ dans $\varphi_K A$.

Proposition 2.5.- Si $K \supset K^2 = Q$ est un sous-demi-groupe, chacune des propriétés suivantes est simultanément vraie ou non pour $Q \subset A$ et pour $\psi_Q : Q \subset \psi_Q A$: Q est libérable, unitaire, net, absorbant.

Soient $B \supset R$ deux demi-groupes quelconques et ψ un homomorphisme. Une condition nécessaire et suffisante pour que quel que soit $x \in B$, $\psi(Qx^{-1}) = \psi Q(\psi x)^{-1}$, est que Q soit saturé pour l'équivalence d'homomorphisme attachée à ψ . Car :

$a \in Qx^{-1} \Leftrightarrow ax \in Q$ qui entraîne : $\psi a \psi x \in \psi Q$ soit $\psi a \in \psi Q(\psi x)^{-1}$; et réciproquement : $\psi a \psi x \in \psi Q$ entraîne $ax \in Q$ est la définition même de la saturation.

Or les propriétés indiquées ne font appel qu'à l'opération \dots
 $Qx^{-1} = \bigcup_{x \in X} Qx^{-1}$.

Proposition 2.6.- Soient $A \supset K$ et ψ un homomorphisme. Une condition nécessaire et suffisante pour que $\varphi_K = \varphi_{\psi K} \circ \psi$ au sens de la composition (\circ) des homomorphismes est que K soit saturé pour l'équivalence σ attachée à ψ (ce que l'on peut noter $\psi^{-1} \psi K = K$).

La condition est nécessaire : si $a \in A - K$, $k \in K$ et $\psi k = \psi a = k' \in \psi K$ alors $(\varphi_{\psi K} \circ \psi) k = (\varphi_{\psi K} \circ \psi) a$ et $\varphi_K a \neq \varphi_K k$.

La condition est suffisante : si σ est moins forte que $\equiv (A, K)$ $a \equiv b (A, K)$ entraîne $\psi a \equiv \psi b (\psi A, \psi K)$.

Définition. Si A est un demi-groupe libre et K un complexe de A nous dirons que $A \supset K$ est une extension libre de $A' \supset K'$ s'il existe un homomorphisme ψ tel que (1) $A' = \psi A$, (2) $K' = \psi K$, (3) $\psi^{-1} \psi K = K$.

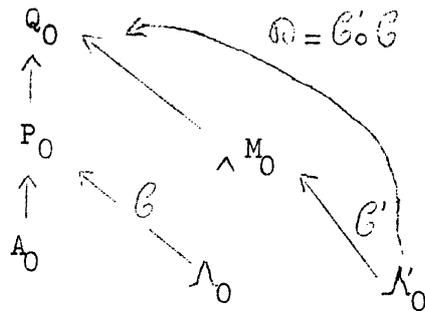
Proposition 2.7.- L'ensemble des couples $(A' \supset P')$ tels que leurs GSF soient isomorphes au couple syntaxiquement simple $(\bar{A} \supset \bar{P})$ est identique à l'ensemble des couples $(\psi A \supset \psi P)$ où $(A \supset P)$ est une extension libre de $\bar{A} \supset \bar{P}$ et où $\psi^{-1} \psi P = P$.

En effet d'après 2.6 on a : d'une part $\varphi_P A = \bar{A}$ et $\varphi_P P = \bar{P}$ d'autre part $\varphi_P = \varphi_{\psi P} \circ \psi$.

Remarque 1. 2.7 montre que l'on peut construire pour toute puissance α de l'alphabet assez grande au moins un code admettant comme GSF un couple syntaxiquement simple $\bar{A} \supset \bar{P}$ où \bar{P} est libérable dans A donné :

Si $\bar{g}_1 \dots \bar{g}_m$ sont des générateurs de \bar{A} et \bar{h}_j d'autres éléments de \bar{A} , il suffit de faire correspondre les lettres $a_{1_1} a_{1_2} \dots a_{1_i}$ à \bar{g}_1 les lettres $a_{2_1} a_{2_2} \dots a_{2_i}$ à $\bar{g}_2 \dots$, $a_{m_1} a_{m_2} \dots a_{m_i}$ à \bar{g}_m et $a_{j_1} \dots a_{j_i}$ à \bar{h}_j en posant $\Psi a_{\ell,k} = \bar{a}_\ell \in \bar{A}$ et de définir P par $\Psi^{-1} \bar{P}$ ce qui assure que $\Psi^{-1} \Psi P = P$. On observera qu'en général pour un tel choix quelconque de générateurs l'ensemble des mots $P_0 = P - ((P - \phi)^2 \setminus \phi)$ n'est pas borné même si \bar{A} est fini.

Remarque 2. Un problème important est celui du surcodage : c'est-à-dire celui de l'utilisation des mots d'un code primaire comme "lettres" d'un alphabet secondaire ainsi qu'il est indiqué dans le schéma ci-contre, qui est illustré



par les identifications suivantes :

A_0 : l'"alphabet" morse formé des trois

"lettres primaires" : point, trait
intervalle

Λ_0 : l'ensemble des signes typographiques mis en correspondance par G avec :

P_0 : l'ensemble des séquences de points traits et intervalles qui

représentent des signes typographiques dans le code Morse. P_0 est à la fois un ensemble de mots primaires et un alphabet secondaire.

Λ'_0 : l'ensemble des mots du français écrit mis en correspondance

1°) par G' avec M_0 : l'ensemble des séquences de signes typographiques qui représentent des mots du français écrit (M_0 est un ensemble de mots primaires pour G').

2°) par $G_0 = G \cdot G$ avec Q_0 : l'ensemble des suites de points, traits intervalles, qui représentent des mots du français écrit (Q_0 est un ensemble de mots secondaires par rapport à P_0 et primaires par rapport à \mathcal{R}).

Il serait utile de caractériser au moins partiellement les GSF de $(A \supset Q)$ au moyen de ceux des "facteurs" $(A \supset P)$ et $(\Lambda \supset M)$. Le problème sera étudié dans un autre travail et les énoncés suivants très simples sont avec 2.6 et 2.7 à la base des démonstrations.

Proposition 2.8. - Si K et K' sont deux complexes de A $a \geq b$ (A, K) et $a \geq b$ (A, K') entraînent $a \geq b$ ($A, K \cap K'$)

Immédiat : car $xy \in K \cap K'$ entraîne $xy \in K$ et $xy \in K'$ qui entraînent $xay \in K$ et $xay \in K'$ donc $xay \in K \cap K'$.

On observera que même si K et K' étaient des demi-groupes P et P' il serait possible que $P \supset P'$ sans que pour autant $\cong (A, P)$ soit plus forte que $\cong (A, P')$ au contraire de ce qui se produit si A est un groupe. En effet, l'intersection de $\cong (A, P)$ et de $\cong (A, P')$ si $(P' \subset P)$ est identique à l'intersection de $\cong (A, P')$ et $\cong (A, P - P')$ qui ne sont pas en général identiques.

Proposition 2.9. - Si $A' \subset A$ et $K' = A' \cap K$ la restriction de $\cong (A, K)$ à A' est plus forte que $\cong (A', K')$. En effet, $xy \in K'$ implique $xay \in K$ par hypothèse et, si $a, b, x, y \in A'$, on a donc $xay \in A \cap K = K'$.

3.- Préfixes.

Définition. On appellera "préfixe à droite" ($\Pi_i^* \in \hat{\Pi}^*$) (respectivement préfixe à gauche : ${}^* \hat{\Pi}_j \in {}^* \hat{\Pi}$) les classes d'équivalence de A pour la relation $\sim^* (A, P)$ (respectivement ${}^* \sim (A, P)$) définie par : $a \sim^* b (A, P)$ si et seulement si $a^{-1}P = b^{-1}P$ (respectivement : $a^* \sim b (A, P)$ si et seulement si $Pa^{-1} = Pb^{-1}$).

Il est classique [3] que :

Proposition 3.1. - La relation $\sim^* (AP)$ est régulière à droite (respectivement : ${}^* \sim (AP)$ est régulière à gauche) et $\cong (A, P)$ est plus forte que l'intersection de $\sim^* (AP)$ et ${}^* \sim (A, P)$.

En effet $a \cong b (A, P)$ s'écrit aussi bien :

$$(ax)^{-1}P = (bx)^{-1}P \text{ pour tout } x \text{ que } P(xa)^{-1} = P(xb)^{-1} \text{ pour tout } x.$$

Il en résulte :

Proposition 3.2. - La représentation A^* des éléments $x \in A$ comme application de l'ensemble $\hat{\Pi}^*$ des préfixes à droite (à gauche) dans lui-même est une représentation isomorphe de \bar{A} .

En effet puisque \sim^* est régulière à droite $a, b \in \hat{\Pi}_i^*$ entraîne $ax, bx \in \hat{\Pi}_j^* = \hat{\Pi}_i^* x$ pour un certain J quels que soient a, b et x , d'autre part $\hat{\Pi}_i^* x = \hat{\Pi}_i^* y$ pour tout $\hat{\Pi}_i^* \in \hat{\Pi}$ entraîne $x \cong y (A, P)$ et réciproquement.

On déduit de cette remarque la :

Proposition 3.3.- $A \supset P$ étant deux demi-groupes quelconques si $\psi_P P = \bar{P}$ est un groupe $A = \psi_P A$ s'il est fini, est la réunion d'un groupe et éventuellement d'un zéro.⁽⁶⁾

En effet la représentation 3.2 étant isomorphe \bar{P} est un groupe si et seulement si toutes les applications correspondantes de $\hat{\pi}^*$ dans lui-même sont des permutations. Donc si $a \in A$ est tel que $xay \in P$ pour au moins un couple $x, y \in A$, a est aussi une permutation. Sinon $ab \equiv ba \equiv a$ (A, P) pour tout b et $\psi_P a = 0$.

Remarque 1. Il est possible et souvent commode de développer la théorie de l'équivalence syntaxique de la façon suivante : soit $a \in A$. L'ensemble $a^{-1}P$ est celui des séquences $b \in A$ telles que $ab \in P$ et l'ensemble $\hat{\pi}_a^* = P(a^{-1}P)^{[-1]}$ celui des séquences c telles que $cb \in P$ pour tout $b \in aP$.

Evidemment $a \in \hat{\pi}_a^*$ et $\hat{\pi}_a^{[-1]}P = a^{-1}P$. On peut vérifier que $\hat{\pi}_a^*$ ainsi défini est précisément le préfixe à droite qui contient a . De la même manière on trouverait ${}^*\hat{\pi}_a = (Pa^{-1})^{[-1]}P$, et la relation ρ entre préfixes à droite et à gauche définie par : $\hat{\pi}_i^* \rho {}^*\hat{\pi}_j$ si et seulement si $a \in \hat{\pi}_i^*$ et $b \in {}^*\hat{\pi}_j$ entraîne $ab \in P$, permet d'établir une correspondance de Galois [2] entre $\hat{\pi}^*$ et ${}^*\hat{\pi}$. Le treillis complet associé pourrait être appelé le "treillis fondamental" du code $(A \supset P)$ et ses propriétés, une fois encore, ne dépendent que des GSF $(\bar{A} \supset \bar{P})$.

Remarque 2. Il est utile de distinguer certains préfixes remarquables :

$\hat{\pi}_0^*$: le résidu de P dans A au sens de F. Dubreil [3] : l'ensemble des a tels que ax n'appartienne à P pour aucun x .

$\hat{\pi}_\infty^*$: le préfixe "absorbant" : l'ensemble des a tels que $ax \in P$ pour tout x .

$\hat{\pi}_1^*$: le préfixe unité tel que $a \in \hat{\pi}_1^*$ entraîne $a \in P$ et $ax \in P$ seulement si $x \in P$.

$\hat{\pi}_{P_i}^*$: les préfixes qui sont des sous-classes de P ($a \in \hat{\pi}_{P_i}^*$ entraîne $a \in P$).

L'existence ou la non-existence de tels préfixes non vides est trivialement liée au fait que le code est net, absorbant, unitaire etc.

⁽⁶⁾ un zéro est un élément 0 tel que $x0 = 0x = 0$ pour tout x .

Exemple. Soit le code suivant (unitaire à droite) :

$$C_\alpha = a ; \quad C_\beta = ab ; \quad C_\gamma = bb \\ \pi_0^* \ni b ; \quad \pi_\infty^* \ni a ; \quad \pi_1^* \ni bb ; \quad \pi_{p\beta} \ni ab .$$

Proposition 3.4. - Si l'alphabet d'un code est fini, et si la longueur de ses mots est bornée par $L < \infty$ alors son GSF \bar{A} est fini.

Puisque \bar{A} possède une représentation comme groupe d'application de $\bar{\pi}$ (7) dans lui-même, il suffit de montrer que $\bar{\pi}$ est fini. Or si $ax \in P$, ou bien la longueur de x est $\leq L$, ou bien il existe x' diviseurs à gauche de x de longueur inférieure à L tels que $ax' \in P$. Donc $a^{-1}P = b^{-1}P$ si et seulement si $a^{-1}P \cap X_L = b^{-1}P \cap X_L$ où X_L est l'ensemble fini des séquences de lettres de longueur $\leq L$.

Proposition 3.5. - Dans tout code unitaire l'ensemble des préfixes différents du résidu est une image homomorphe de l'ensemble des diviseurs (à gauche) des mots.

Montrons d'abord que $p \in P$ entraîne $(pa)^{-1}P = a^{-1}P$ pour tout a . Or $pax \in P$ entraîne $ax \in P$ par hypothèse puisque P est unitaire.

Soit donc $a \in \pi_i \neq \pi_0$; il existe y tel que $ay \in P$ donc a est diviseur à gauche d'une suite de mots donc, ou bien a est lui-même diviseur à gauche d'un mot, ou bien il existe $p \in P$, tel que $a = pa'$ et $a' \in \pi_i$, d'où le résultat par récurrence. On en déduit la remarque utile suivante :

Proposition 3.6. - Il existe une correspondance biunivoque entre les codes unitaires et les structures d'arbres enracinés (8), les codes nets, correspondant

(7) Pour simplifier et puisque dans la pratique un ordre temporel est toujours prescrit pour la suite des lettres, nous conviendrons dorénavant sauf mention expresse du contraire, que préfixe signifie, préfixe à droite; unitaire : unitaire à gauche; net : net à droite. Cette convention est la plus simple quand l'ordre temporel est identifié avec l'ordre gauche \rightarrow droite.

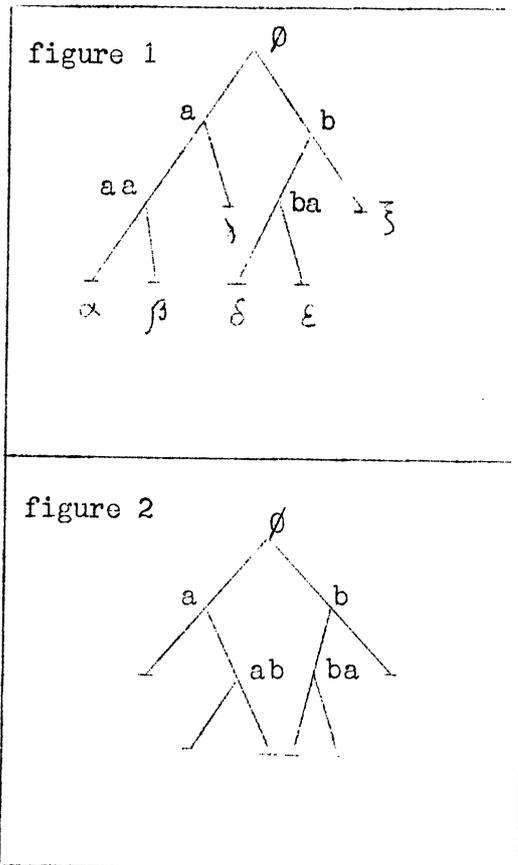
(8) Soit U un ensemble ordonné par ρ . Pour tout u soit $\rho[u] =$ l'ensemble des $u' \in U$ tels que $u' \rho u$. U est un arbre enraciné si 1°) $\bigcap_{u \in U} \rho[u] = u_1 \neq \emptyset$ (u_1 est la racine) 2°) pour tout $\rho[u]$ est un ensemble totallement ordonné. Si $u \in \rho[v]$ entraîne $v = u$ u est une extrémité, sinon u est un noeud. Si pour tout u qui n'est pas une extrémité il existe un nombre fixe $K < \infty$ de v tels que 1°) $u \rho v$ et 2°) $u \rho w \rho v$ entraîne $w = u$ ou $w = v$, l'arbre est dit complet.

aux arbres complets.

Il suffit d'identifier la racine à la séquence vide, les noeuds aux diviseurs à gauche des mots, et les extrémités aux mots eux-mêmes. On observera que deux diviseurs appartiennent au même préfixe si et seulement si les sous-arbres correspondant dont ils sont les racines sont identiques.

Exemple : soit le code suivant :

$C_\alpha = aaa$; $C_\beta = aab$; $C_\gamma = ab$; $C_\delta = baa$; $C_\epsilon = bab$; $C_\zeta = bb$
 il lui correspond l'arbre ci-contre (figure 1). Les diviseurs "a" et "b"



appartiennent au même préfixe.

On prendra garde que des arbres équivalents à des permutations près de noeuds correspondent en général à des codes dont les GSF sont différents. (l'arbre de la figure 2 par exemple est équivalent par permutation à celui de la figure 1 mais son code à un autre GSF) .

Une conséquence immédiate de la représentation des codes unitaires par des arbres enracinés est :

Proposition 3.7.- Dans un code unitaire une condition nécessaire et suffisante pour que tous les mots aient une longueur bornée est que, si a et ax appartiennent au même préfixe, il existe au moins un diviseur à gauche x' de x tel que $ax' \in P$.

En effet les sous-arbres a et ax ne sont jamais identiques, si $ax' \notin P$, pour tous les diviseurs à gauches de x' .

Proposition 3.8.- Une condition nécessaire et suffisante pour qu'un code unitaire dont la longueur des mots est bornée par $L < \infty$ possède une suite non vide de lettres ℓ telle que φ_P^ℓ soit un élément neutre bilatère de φ_P^A est qu'il existe une lettre $a \in A_0$ telle que les séquences $a, a^2 = aa, a^3 = aaa, \dots, C = a^m \in P$ forment un système de représentants de ses préfixes, à l'exception éventuelle du résidu π_0 .

Par hypothèse ℓ est une application de $\overline{\Pi}$ sur lui-même, si donc ℓ contient la lettre a dans son expression, a est une permutation des préfixes (Cf. 3.3) et puisque la longueur des mots est bornée il existe un m minimum fini tel que $a^m \in P$ et $a^m \equiv \varrho(A, P)$. On vient de voir que $a^{m'} \sim a^{m''}$, $m', m'' \leq m$ entraîne $m' = m''$. Soit $b \in \overline{\Pi}'$. Pour au moins un $n' < \infty$, $ba^{n'} \in P$ mais aussi pour un $m' = m - n'$ (à un multiple de m près) $a^{m'} a^{n'} \in P$ donc $\overline{\Pi}'$ est identique au préfixe contenant $a^{m'}$ puisque la multiplication à droite par $a^{n'}$ est une permutation. On en déduit :

Proposition 3.9. - Si la longueur des mots d'un code est bornée et si la puissance de son alphabet est finie, une condition nécessaire et suffisante pour que son GSF \overline{A} soit un groupe est que $P_0 = X\varrho$ (= l'ensemble de toutes les séquences de longueur ℓ : code uniforme de longueur ℓ). Dans ce cas \overline{A} est le groupe cyclique d'ordre ℓ et P se réduit à son élément neutre.

De 3.7, il résulte que toutes les lettres de l'alphabet doivent correspondre à des permutations circulaires de même ordre ℓ sur les préfixes. Montrons que $a^n b \sim a^{n+1}$ quels que soient les lettres a et b , (il sera convenu que les indices seront des entiers positifs réduits modulo ℓ).

Supposons que, pour un n' et un m' , $m' \leq n'$, on ait $a^{n'} b \sim a^{m'}$ et considérons les séquences non bornées $c_x = a^{n'} (ba^{m'-n'})^x$ (n' fois suivis de x fois la sous-séquence b suivi de $m' - n'$ a).

On vérifie facilement qu'aucun de leurs diviseurs à gauche n'appartient à P (ils sont tous \sim à un $a^{n''}$ où $n'' \leq n'$). Donc $c_x a^{\ell-n'} \in P$ est un mot contrairement à l'hypothèse selon laquelle la longueur des mots est bornée. Donc pour tout n' $a^{n'} b \sim a^{m'}$ où 1°) m' parcourt avec n' l'ensemble de tous les nombres $\leq \ell$ (car b est une permutation !) 2°) $m' > n'$ si $n' \neq \ell$. Ceci n'est possible que si $m' = n' + 1$ (modulo ℓ) ou encore que si $a \equiv b (A, P)$ puisque $\overline{\Pi}_i a \sim \overline{\Pi}_i b$ pour tout préfixe. Donc \overline{A} est un groupe fini à un seul générateur, etc.

Remarque.

Le résultat précédent pose le problème de savoir s'il existe des couples syntaxiquement simples $\overline{A} \supset \overline{P}$ où \overline{P} serait unitaire et net à droite et à gauche sans que \overline{A} soit un groupe et tels que pour au moins un choix de générateurs la longueur des mots soit bornée. La proposition 3.10 est une réponse

très partielle à cette question. Nous montrons d'abord

Proposition 3.9.- La condition nécessaire et suffisante pour que l'élément u du demi-groupe libre A satisfasse pour au moins un couple $v, w \neq \emptyset$ à : $a = uv = wu$ où la longueur ℓ de a est inférieure ou égale au double de la longueur ℓ' de u est que u soit de la forme $x(yx)^n$.

Raisonnons par récurrence : les longueurs de v et w étant strictement inférieures à $\ell/2$ on a :

$$a = wrv \quad \text{pour un certain } r \quad \text{et} \quad u = wr = rv.$$

Si la longueur de r est plus petite ou égale $\ell'/2$ on a encore : $u = rsr$ avec $w = rs$ et $v = sr$. Sinon on est ramené au problème initial avec u au lieu de a et r au lieu de u . Comme les longueurs des séquences ne peuvent que décroître strictement on a bien le résultat.

Proposition 3.10.- Pour un alphabet de $K < \infty$ lettres il existe au moins un code du type qui vient d'être décrit et dont le nombre de mots est égal à $1 + K^\ell - 2K^{\ell''} + K^{2\ell''}$ pour tout $\ell \geq 3$ et $\ell'' < \ell/2$.

Nous considérons d'abord le code uniforme \mathcal{U} dont P_0 est identique à l'ensemble X_ℓ de toutes les séquences de ℓ lettres. Il satisfait aux conditions voulues sauf que son GSF est un groupe.

Soit u une séquence fixe de longueur $\ell' < \ell$ et soient les ensembles suivants :

P'_0 : les mots (de $X_\ell = P_0$) dont u est un diviseur à gauche ($P'_0 = u(u^{-1}P_0)$)

P''_0 : les mots dont u est un diviseur à droite mais non à gauche
 $(P''_0 = (P_0 u^{-1})u - P'_0 \cap (P_0 u^{-1})u)$

X : l'ensemble des séquences de longueur $\ell'' = \ell - \ell'$

Q'_0 : l'ensemble $P''_0 X_{\ell''}$ des séquences de la forme $q = pv$ avec $p \in P''_0$ et $v \in X_{\ell''}$.

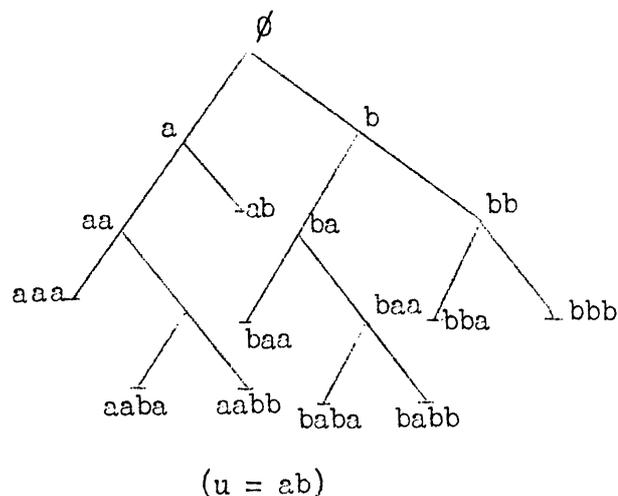
Considérons le code \mathcal{B} dont l'ensemble Q_0 des mots est la réunion de u , Q'_0 et de l'ensemble $Q''_0 = P_0 - (P'_0 \cup P''_0)$. Par construction \mathcal{B} est unitaire à gauche et net à droite. D'autre part si $q, q' \in Q_0$ q ne peut évidemment pas être un diviseur à droite de q' dans les cas suivants :

1°) $q' = u$; 2°) $q \in Q'_0$; 3°) $q \in Q''_0$ et $q' \in Q''_0$; 4°) $q = u$ et $q' \in Q''_0$

Il reste à discuter 5°) le cas de $q \in Q_0''$ et $q' \in Q_0'$ mais les mots de Q_0'' sont de la forme $xuv = xp$ avec $v \in X_{\ell''}$ et $p \in P_0'$ et $p \in P_0'$ ne peut admettre aucun diviseur propre à droite de longueur ℓ . 6°) \mathcal{B} sera donc unitaire à droite si et seulement si $Q_0' u^{-1} = \emptyset$ ou encore si $uvu^{-1} = \emptyset$ pour tout $v \in X_{\ell''}$. Ceci n'est possible comme on l'a vu en 3.9 que si $\ell' > \ell/2$ et si u n'est pas de la forme exceptionnelle $x(yx)^n$. Le calcul du nombre des mots n'offre aucune difficulté.

Exemple.

Le cas le plus simple ($K = 2$; $\ell = 3$) est décrit par l'arbre suivant. Il contient 9 mots et son GSF (privé de l'élément neutre bilatère) a 24 éléments et ne possède pas d'idéaux propres. Cette dernière particularité n'est pas une nécessité pour les GSF des codes de ce type.



4.- Méthodes de dénombrement.

On supposera désormais toujours que l'alphabet a une puissance $K < \infty$. On dira, pour abrégé, qu'un code est borné si la longueur de ses mots est plus petite que $L < \infty$.

Nous commencerons par compléter et systématiser divers résultats plus ou moins explicitement connus et utilisés par les auteurs qui ont étudié ces questions et notamment B. Mandelbrot [9].

Notations $g(t) = 1 - \sum_{i=1}^{\infty} n_i t^i$: (avec n_i = nombre des mots de longueur i)

La fonction de structure du code $G(t) = 1 + \sum_{i=1}^{\infty} N_i t^i$ (avec N_i = nombre des éléments de P de longueur i) : la fonction génératrice du code.

$H(t) = |Dt - 1|$ la fonction caractéristique du code avec D : la somme des matrices correspondant aux lettres de A_0 dans la représentation régulière de \bar{A} .

$h_d(t)$ et $h_g(t)$: les déterminants correspondants dans les représentations de \bar{A} comme demi-groupe d'application dans eux-mêmes de $\bar{\Pi}^*$ et de ${}^*\bar{\Pi}$.

Proposition 4.1. - $g(t)$ est un diviseur commun et $H(t)$ un multiple commun de $h_g(t)$ et $h_d(t)$. On a $G(t) = \frac{1}{g(t)}$. Introduisons encore $N_i(\bar{a})$: nombre de séquences a de longueur i telles que $\varphi_P a = \bar{a} \in P$. Les $N_i(\bar{a})$ sont liés entre eux par un système d'équations aux différences finies :

$$N_{i+1}(\bar{a}) = \sum_{\bar{b} \in \bar{A}} \delta_{\bar{a}, \bar{b}} N_i(\bar{a})$$

où $\delta_{\bar{a}, \bar{b}} = 0$ ou $1, 2, \dots, K$ selon qu'il existe $0, 1, 2, \dots, K$ lettres a_i dans A_0 telles que $\varphi(a_i) = \bar{b} \in \bar{A}$. La matrice des $\delta_{\bar{a}, \bar{b}}$ est précisément D et si les ρ_J sont les racines du déterminant $|Dt - 1|$ on sait que $N_i(\bar{a}) = \sum_J \beta_{\bar{a}, J} \rho_J^{-i}$ pour un certain système de constantes $\beta_{\bar{a}, J}$. D'autre part, les $N_i(\bar{\Pi}^*)$ ou les N_i qui sont obtenus comme sommes (par rapport à \bar{a}) de certains $N_i(\bar{a})$ peuvent être calculés directement à partir de représentation de \bar{A} sur $\bar{\Pi}^*$ ou sur ${}^*\bar{\Pi}$ ce qui établit les relations de divisibilité indiquées. Enfin on a directement $N = \sum_i N_{\ell-i} n_i$ ce qui donne $G(t) = 1/g(t)$ par un raisonnement classique. (Cf. 6)

Proposition 4.2. - La fonction de structure d'un code borné unitaire (à gauche) et net (à droite) (code UND) est de la forme :

$$g(t) = (1 - Kt) \bar{g}(t) \quad \text{où : } \bar{g}(t) = 1 + \sum \bar{n}_i t^i \quad \text{avec } 0 \leq \bar{n}_{i+1} \leq K \bar{n}_i.$$

K^{-1} en est donc la plus petite racine positive. Réciproquement tout polynôme $g(t)$ de la forme précédente peut être considéré comme la fonction de structure d'au moins un code UND.

Soit \bar{n}_ℓ le nombre des séquences de longueur ℓ qui sont diviseur propre à gauche des mots d'un code unitaire (= qui correspondent aux "noeuds" de l'arbre de codage). On a : $\bar{n}_0 = 1$ et $\bar{n}_{\ell+1} + n_{\ell+1} \leq K \bar{n}_\ell$ le signe \leq n'étant partout remplacé par le signe = que si l'arbre est complet, c'est-à-dire si le code est net. Multipliant ces égalités par $t^{\ell+1}$ et sommant, on obtient bien : $g(t) = (1 - Kt) \bar{g}(t)$.

Réciproquement, un polynôme de la forme $g(t) = 1 - \sum_{i=1}^{\ell} n_i t^i$ ($n_i \geq 0$) admet au moins un diviseur de la forme $(1 - K't)$ et l'on a :

$g(t) = (1 - K't)(1 + \bar{n}_1 t + \bar{n}_2 t^2 \dots)$ avec $\bar{n}_{\ell+1} \leq K' \bar{n}_i$ pour tout i . Donc si $g(t)$ est un polynôme $0 \leq \bar{n}_{i+1} \leq K' \bar{n}_i$.

D'autre part K' ne saurait être $> K$ car sinon pour ℓ assez grand $N_{\ell} = \sum \beta_j r_j^{-\ell} \geq \beta_1 K'^{\ell}$ serait plus grand que K^{ℓ} , ce qui est impossible puisqu'il s'agit d'un code et que ceci signifierait que à au moins l'une des K^{ℓ} suites de ℓ lettres correspondent plusieurs séquences de mots. Enfin si \bar{n}_i satisfait à $0 \leq \bar{n}_{i+1} \leq K' \bar{n}_i$ pour tout i il est facile de construire un arbre dont \bar{n}_i et n_i sont respectivement les nombres de noeuds et d'extrémités à distance i de la racine ⁽⁹⁾. On observera :

Proposition 4.3. - Dans un code UND le coefficient β_1 de K^i dans l'expression $N_i = \sum \beta_j r_j^{-i}$ est précisément l'inverse de la longueur moyenne \bar{L} des mots quand la fréquence des mots de longueur ℓ est proportionnelle à $K^{-\ell}$.

En effet : $G(t) = \frac{1}{g(t)} = \frac{(\bar{g}(K^{-1}))^{-1}}{1 - Kt} + \frac{B}{\bar{g}(t)}$ où B est une certaine constante.

D'autre part, on a, par hypothèse :

$L = \sum n_i i K^{-i}$ c'est-à-dire que $-L$ est la valeur de $\frac{\partial g(t)}{\partial t} K^{-1}$ pour $t = K^{-1}$ soit encore précisément $\bar{g}(K^{-1})$.

Remarque.

Attachons à tout $a_i \in A_0$ une probabilité $\omega_i \geq 0$ ($\sum \omega_i = 1$) et considérons au lieu de D la matrice $\sum \omega_i D_i$ en appelant D_i la matrice associée à a_i dans la représentation régulière de \bar{A} . Les mêmes considérations sont encore valables et les fonctions $N_i(\bar{a})$ deviennent des probabilités dans un processus stochastique (sur A) ou les lettres successives sont tirées au sort indépendamment et avec les probabilités constantes ω_i . Le cas traité ici correspond à un changement d'échelle sur t (qui devrait être remplacé par $t_1 = tK^{-1}$) à celui où tous les ω_i sont égaux et on vérifiera que les propriétés 4.1, 4.2, 4.3 sont des théorèmes bien connus pour les chaînes de Markoff. Une théorie peut aussi être développée dans laquelle les matrices D_i n'ont pas nécessairement leurs éléments égaux à zéro ou à un, mais il peut alors

⁽⁹⁾ Ceci complète une démonstration de B. Mandelbrot [8] qui laissait ouverte la possibilité que β_1 (le coefficient du terme correspondant à la plus petite racine) soit plus grand pour certains codes non unitaires que pour tout code unitaire.

se produire que la relation d'équivalence $\cong (A, P)$ doive être remplacée par une relation plus faible pour tenir compte de la dépendance en chaîne des lettres.

Nous n'avons besoin ici que des méthodes énumératives pour établir les deux résultats suivants qui relèvent strictement de la théorie algébrique des demi-groupes et qu'il semblerait difficile d'obtenir autrement.

Proposition 4.4.- Tout code borné d'alphabet fini net à droite est unitaire à gauche.

Soit $L < \infty$ la longueur maximum des mots du code. Puisque le code est net, il correspond, à chacune des K^ℓ séquences de longueurs ℓ , au moins une séquence de longueur $\leq \ell + L$ qui appartient à P est dont elle est un diviseur à gauche. Donc : $N_\ell + N_{\ell+1} + \dots + N_{\ell+L} \geq K^\ell$ pour tout ℓ assez grand. Donc 4.1, au moins une des racines de $g(t)$ est $\leq K^{-1}$ et comme il s'agit d'un code cette racine de module minimum est précisément K^{-1} . Donc (4.2), $g(t)$ est identique à la fonction de structure d'un certain code UND.

Considérons maintenant P'_0 , le sous-ensemble des mots du code qui n'admettent aucun autre mot comme diviseur à gauche. Puisque le code est net toute séquence

soit admet un $p \in P'_0$ comme diviseur à gauche

soit est un diviseur à gauche d'un $p \in P'_0$.

Donc l'arbre de codage correspondant au code UND est complet et sa fonction de structure $g'(t)$ admet la racine K^{-1} . Il est impossible que $P'_0 \neq P_0$ car ceci impliquerait que $g(t) = g'(t) - \sum n_i'' t^i$ (où les n_i'' correspondent aux mots de $P_0 - P'_0$) ce qui ne se peut puisque $g(K^{-1}) = g'(K^{-1}) = 0$. Donc $P'_0 = P_0$ et le code est bien UND.

Proposition 4.5.- Si un code est UND

ou bien il est aussi net à gauche (et par conséquent unitaire à droite)

ou bien le code opposé est tel qu'il existe des séquences non bornées dont le premier mot ne peut pas être décodé sans que soit connue la totalité du message.

On a vu (proposition 1.9) que s'il existait une borne L à de telles séquences on aurait $P_n = \emptyset$ pour un $n < \infty$ ou encore (en posant $P_0^m =$ l'ensemble des séquences formées de m mots non vides) que pour un certain n' $P_0^{n'}$ ne contiendrait aucune paire d'éléments dont l'un soit diviseur à gauche de l'autre.

Or il est clair que si P_0 est l'ensemble des mots d'un code net il en est de même de P_0^n et réciproquement car la fonction de structure $g_n(t)$ de P_0^n est donnée par : $1 - g_n(t) = (1 - g(t))^n$ quel que soit P_0 et n'admet la racine K^{-1} que si $g(t)$ l'admet elle-même. Si donc $\{P_0^m\}$ était unitaire, il serait net puisque la fonction génératrice d'un code est la même que celle du code opposé. Donc $\{P_0\}$ lui-même serait à la fois unitaire et net des deux côtés.

BIBLIOGRAPHIE

- [1] - J. RIGUET : Sur la représentation des syntaxes. (à paraître dans : Zeitschrift für math. Log., 1955).
- [2] - J. RIGUET : Relations binaires (Bull. Soc. math. France, t. 76, 1948, p. 114-155).
- P. DUBREIL : Trois mémoires intitulés "Contribution à la théorie des demi-groupes. (I, II, III).
- [3] - I : Mem. Acad. Sci. Paris, t. 63, 1941, p. 1-52.
- [4] - II : Rendiconti di Matematica e delle sue applicazioni, Università di Roma, serie V, vol. X, 1951, p. 183-200.
- [5] - III : Bull. Soc. math. France, t. 81, 1953, p. 289-306.
- [6] - W. FELLER : Introduction to Probability theory. (Wiley, New York 1950) chapitre 12.
- [7] - B. MANDELBROT : Thèse, Paris 1952.
- [8] - B. MANDELBROT : On recurrent noise limiting codes (Proc. Symp. Inf. Networks, Brooklin, 1954, p. 206-221).
- [9] - B. MANDELBROT : Mémoire dans : Proc. Symp. Inf. Theory, 1955 (à paraître chez Butterworth, London).
- [10] - P. ELIAS : Idem
- [11] - D. SLEPIAN : Idem
- [12] - A.A. SARDINAS et G.W. PATTERSON : A necessary and sufficient condition for unique decomposition ... of coded messages. (Univ. Pennsylvania, Res. Div. reports 50.27, march 1950).
-