



ÉCOLE NATIONALE POLYTECHNIQUE

DEPARTEMENT D'ELECTRONIQUE

PROJET DE FIN D'ETUDES

S U J E T

Application de la classification
à la reconnaissance des phonèmes

Proposé par :

M^r. B. Bouseksou

Étudié par :

R. Belkessa
D. Hamoutene

Dirigé par :

M^r. B. Bouseksou

PROMOTION : JUIN 1985

DEDICACES

à toute ma famille

R A B A H

à mes parents

à tous ceux qui me sont chers

D J A M E L

R E M E R C I M E N T S

Nous tenons remercier vivement notre promoteur Mr BOUSEKSOU
POUR toute l'aide qu'il nous a apportée durant la réalisation
de ce projet.

Nous tenons à exprimer aussi notre reconnaissance à Mr AMENAS
DE L.I.N.I pour son précieux concours.

Nous tenons à remercier également Mr FERRAH MENOUEUR pour les
travaux de dactylographie.

S O M M A I R E

INTRODUCTION	PAGE	I
Chapitre I	Etude de la parole	
1-	Introduction	4
2-	Production de la parole	4
3-	Décomposition de la parole	6
Chapitre II	Analyse de la parole en vue de sareconnaissance	
1-	Introduction	10
2-	Analyse spectrale	11
3-	" " par la prédiction lineaire	16
4-	" " cepstrale	21
5-	" " temporelle	30
Chapitre III	La classification	
1-	Introduction	32
2-	Principes de la classification numérique	32
3-	Etude de la distance	33
4-	L'apprentissage	35
5-	La reconnaissance	39
6-	L'analyse en composantes principales	40
Chapitre IV	Application de la classification à la reconnaissance des phonèmes	
1-	Introduction	43
2-	Nuée dynamique	43
3-	Algorithme des k-moyennes	46
4-	Reconnaissance des phonèmes	54
5-	Conclusion	55
Programme d'apprentissage et de reconnaissance		56
CONCLUSION		61
		64

Table des figures

Schématisation de l'appareil vocal	figure	(1)
Spectre d'un son voisé et d'un son non voisé	"	(2)
Vocodeur à canaux	"	(3)
Vocodeur à formants	"	(4)
Modèle numérique de la production de la parole	"	(5)
Densité spectrale d'un modèle d'ordre 12	"	(6a)
Conversion logarithmique de la densité spectrale	"	(6b)
Obtention des coefficients Mel par un vocodeur à canaux	"	(7)
Bancs de filtres triangulaires	"	(8)

I N T R O D U C T I O N

Depuis le développement des ordinateurs, la communication entre l'Homme et la machine n'est plus une utopie. En effet le mode de communication privilégié entre les hommes qu'est la parole a beaucoup préoccupé les chercheurs.

La communication entre la machine et l'Homme ^{pose} des problèmes très différents suivant le sens dans lequel elle se fait :

- Dans le sens machine-Homme, il s'agit de synthétiser une voix qui soit à la fois intelligible et le moins possible artificielle.

- Dans le sens Homme-machine, il s'agit de permettre à la machine de reconnaître le sémantisme de la parole. Or, la reconnaissance de ce contenu sémantique passe par un certain stades intermédiaires correspondant aux divers niveaux d'information accessibles dans la parole

L'avènement de nouveaux moyens de calcul a incité les chercheurs à les utiliser pour le traitement numérique de la parole. Cet intérêt ne s'est infléchi, bien au contraire, des nouvelles méthodes de traitement du signal en général et du signal de la parole en particulier ne cessent de voir le jour.

Le champ de recherche sur la parole est vaste, il exige des connaissances multidisciplinaires.

Au début, l'intérêt manifesté par les universitaires était purement scientifique, ensuite s'est ajouté l'intérêt économique de dizaines d'entreprises qui ont stimulé les recherches dans ce domaine. En effet, plusieurs entreprises commercialisent actuellement des circuits intégrés spécialisés, des cartes de reconnaissances et de synthèses à usage général et enfin l'intérêt stratégique des militaires.

Dés résultats encourageants ont été obtenus dans le domaine de la synthèse, cela est dû au fait que le mode de production de la parole ainsi que les modes articulatoires combinés de tous les organes du conduit vocal sont assez bien compris.

Dans le domaine de la reconnaissance automatique de la parole, quoique les recherches s'intensifient, l'objectif final n'est pas encore atteint à savoir réaliser un système qui a les caractéristiques suivantes:

- Compréhension de la parole continue
- Vocabulaire de mots (avec syntaxe) d'assez grande taille
- Multilocuteur
- Un taux d'erreur (de sémantique) très faible (moins de 10%)
- Réponse en temps réel

Un tel système n'existe pas actuellement car le processus par lequel l'être humain décode la parole en un ensemble de traits phonétiques n'est pas encore connu.

En effet la perception de la parole n'est pas seulement due à une extraction passive des traits phonétiques, elle met en jeu un phénomène de mémorisation comprenant la génération d'hypothèses de synthèse interne avec comparaison au signal reçu.

La perception de la parole par l'être humain passe par trois phases

- 1- une phase d'acquisition (capteur)
- 2- la paramétrisation réalisée par l'oreille moyenne et l'oreille interne.
- 3- la décision ou la reconnaissance effectuée par le système nerveux central.

Ces trois phases se retrouvent dans tout système de reconnaissance automatique de la parole.

Le problème de la reconnaissance de la parole est en fait un problème de reconnaissance de forme. Il s'agit de trouver une méthode automatique qui partage un ensemble de données en sous ensembles ayant chacun un nom. Notre travail consiste à étudier et à élaborer un algorithme de classification automatique (algorithme des K - Moyennes) pour la reconnaissance des phonèmes. Ce dernier permet après décodage lexical de reconnaître des mots ou des séquences de mots.

CHAPITRE I

ETUDE DE LA PAROLE

I-1)- INTRODUCTION.

" Nous allons parler succinctement, dans ce qui suit ; de l'origine et de la nature de la parole, nous préciserons aussi certaines propriétés du signal de la parole.

1-2--)- PRODUCTION DE LA PAROLE.

1-2 -1)-FONCTIONNEMENT DE L'APPAREIL VOCAL:

La parole, en tant que phénomène physique, résulte de l'excitation du conduit vocal par deux types de sources sonores. Le conduit vocal est une suite de cavités qui servent de résonateurs et suivant leurs formes, il apparait des resonances à des fréquences variables appelées formants, nous distinguons la cavité pharyngale, la cavité buccale et enfin la cavité nasale. Cette dernière est soit inutilisée soit elle se met en dérivation sur la cavité buccale par l'abaissement du vélum. Cela se traduit par l'apparition d'antiformants sur le spectre des signaux (sons nasalisés).

La production des " sons voisés", comme les voyelles, fait intervenir une source d'impulsions périodiques constituée par l'ensemble poumons cordes vocales. Les sons " non voisés", comme certaines consonnes, sont engendrés par une source de bruit, c'est à dire de signaux de forme aléatoire. Le bruit est produit par l'action du courant d'air issu des poumons sur les parois du conduit vocal, soit au niveau d'un resserrement, soit au niveau d'une fermeture totale de celui-ci.

I-2-2) ETUDE SPECTRALE

Les deux sources que nous avons décrites possèdent des caractéristiques spectrales différentes. La source périodique, ou source vocale délivre des oscillations de "relaxation" qui prennent naissance au niveau des cordes vocales.

La première raie du spectre de la source se trouve à la fréquence fondamentale F_0 (pitch) tandis les raies suivantes espacés de F_0 sont des harmoniques. La fréquence fondamentale varie entre 70 et 150 Hz. pour les hommes et de 150 Hz à 300 Hz pour les femmes, et peut dépasser 400 Hz chez les enfants.

.../...

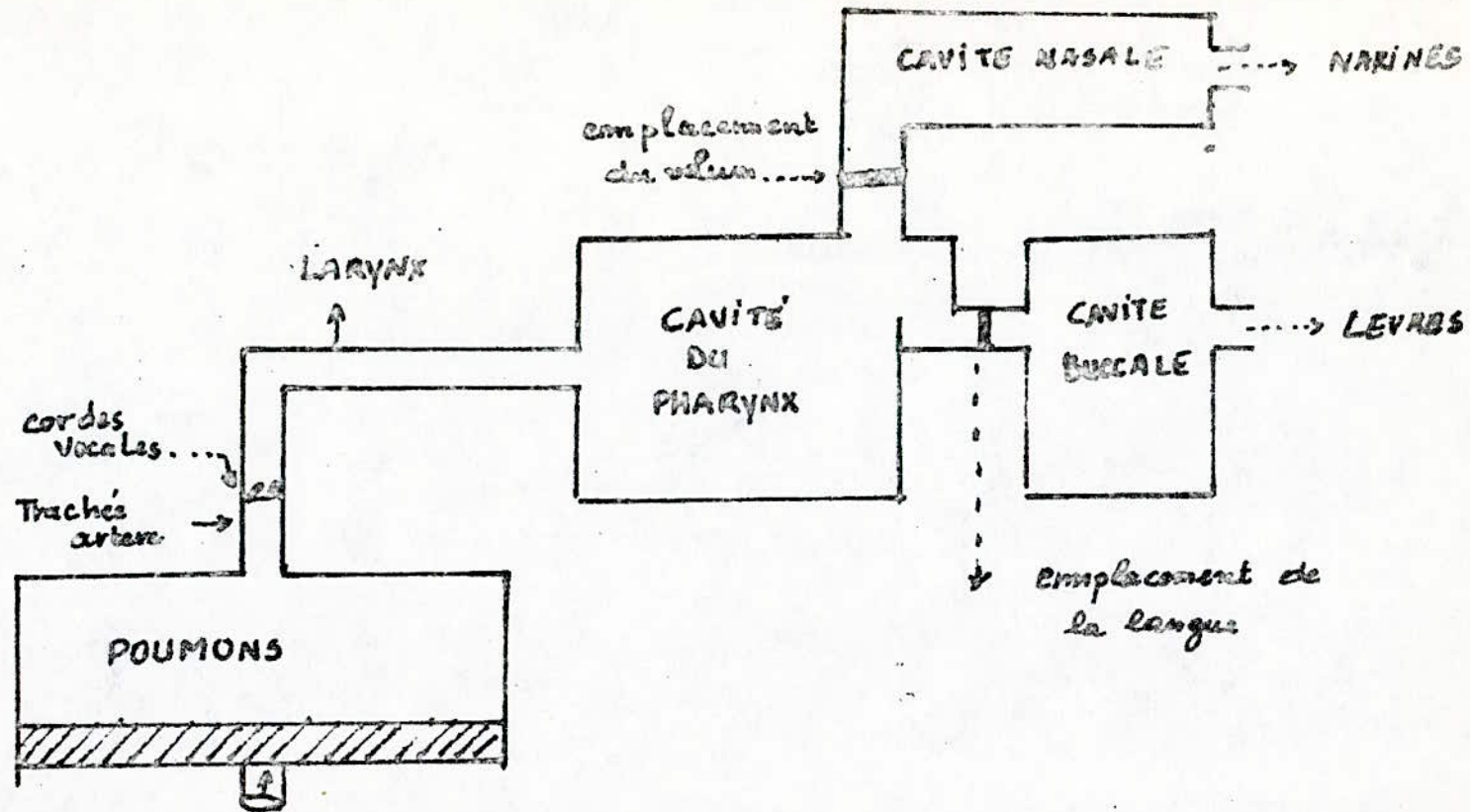


Fig (1) SCHÉMATISATION DE L'APPAREIL VOCAL

La source de bruit est localisée sur le parcours du conduit vocal et sa situation peut varier de 200 à 900 Hz, le second de 500 à 2500 Hz. Les fréquences des formants sont stables lors de la production d'une voyelle et dépendent de cette dernière. En revanche les fréquences des formants changent lors de la production de certaines consonnes car le conduit vocal évolue rapidement dans le temps. Une analyse spectrale de la parole, que l'on peut obtenir à l'aide d'un sonographe, met en évidence les évolutions des formants.

I-3) D E C O M P O S I T I O N D U L A P A R O L E

La phonétique considère depuis longtemps l'existence de sons élémentaires appelés phonèmes. Une trentaine de phonèmes pourrait décrire la prononciation française, les phonèmes du français, au nombre de 36, sont classés en quatre groupes:

-Les voyelles : sons périodiques dont la fréquence du fondamentale définit la hauteur du son. La nature de la voyelle est déterminée par la forme des cavités.

-Les consonnes fricatives: caractérisées par un rétrécissement du passage de l'air. Ce rétrécissement, qui a lieu dans certains endroits du conduit vocal, est équivalent à une source de bruit mise en forme par les cavités de conduit. Si en plus de cette source de bruit les cordes vocales interviennent, la consonne produite est dite voisée. Z/S, V/f, J/ ch.

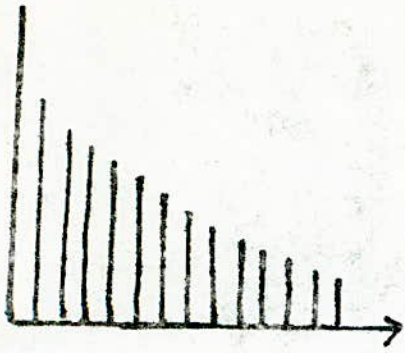
-Les consonnes plosives: sons purement transitoires, ce sont des consonnes momentanées qui supposent une occlusion complète suivie d'une ouverture brusque, semblable à une explosion. Il existe des couples de consonnes plosives semblables en ce qui concerne le caractère voisé ou non b/p, d/t, g/k.

-Les nasales; sons caractérisés par la mise en parallèle des fosses nasales sur le conduit vocal. "In", "an" sont des voyelles nasales, "m" "n" sont des consonnes nasales.

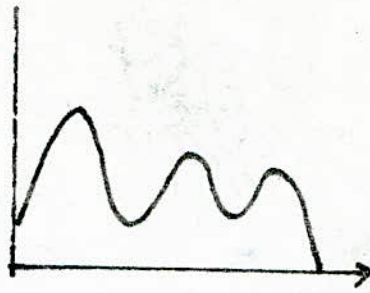
Enfin, il existe quelques cas particuliers : les liquides (l) les vibrantes (V) , les diphtongues.

.../....

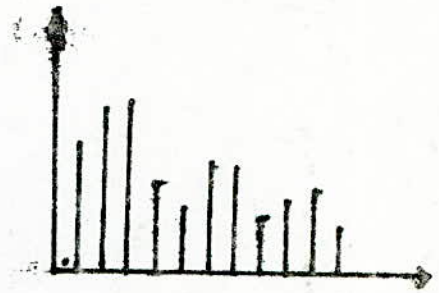
SONS VOISES



SPECTRE DE LA
SOURCE

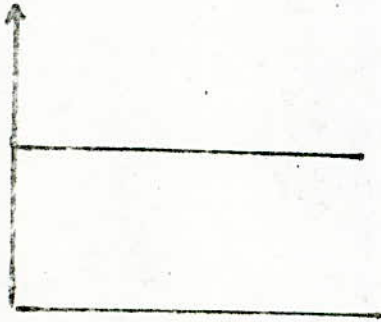


CANAL

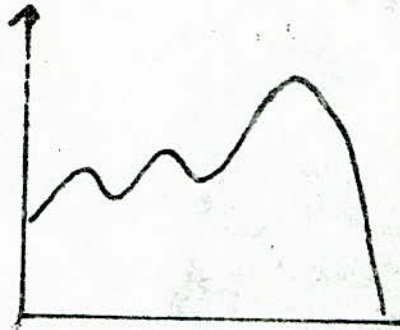


SIGNAL

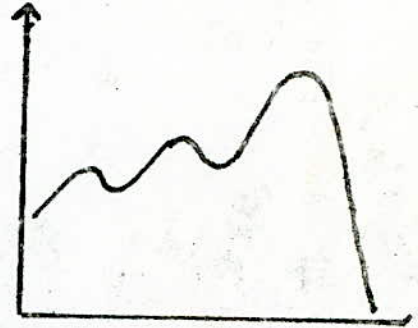
SONS NON VOISES



SPECTRE DE LA
SOURCE



CANAL



SIGNAL

Fig (2)

Pour décrire phonétiquement une phrase parlée, il suffit de reconnaître chacun des phonèmes qui la compose, le plus souvent, les systèmes de reconnaissance de la parole utilise la "segmentation". Cette dernière consiste à découper le signal contenu en unités plus petites qui représenteraient des phonèmes. Au préalable, une segmentation séparant la parole du bruit est nécessaire. Le problème dans ce cas, est de détecter le début et la fin de la phrase. Cette tâche est aisée lorsque le rapport signal sur le bruit est de l'ordre de 60dB mais devient plus difficile lorsque le rapport est inférieur à 30 dB.

La méthode la plus courante, pour la segmentation, est celle qui utilise l'amplitude ou l'énergie du signal.

LES SONS DU FRANCAIS

VOYELLES ORALES			CONSONNES		
PHONETIQUE	EXEMPLE	CODE MACHINE	PHONETIQUE	EXEMPLE	CODE MACHINE
[i]	HABIT	~I	[p]	PERE	-P
[e]	CAFÉ	~E	[t]	TU	-T
[ɛ]	PROCÉS	AI	[k]	CAS	-K
[a]	PAPA	^A	[b]	BON	-B
[ɑ]	VASE	Idem	[d]	DIRE	-D
[ɔ]	ROBE	AU	[g]	GOUT	-G
[o]	POT	'O	[f]	FEU	-F
[u]	MOU	OU	[v]		-V
[y]	TU	~Y	[s]	SORT	-S
[ø]	PEU	EU	[z]	ZERU	-Z
[œ]	PEUR	Idem	[ʒ]	CHOU	CH
[ɚ]	LE	EE	[ʃ]	JOUR	ʒʒ
	VOYELLES NASALES		[ʒ]		-R
[ɛ̃]	PAIN	~N	[l]		-L
[œ̃]	BRUN	Idem	[m]		-M
[ɑ̃]	BLANC	AN	[n]		-N
[ɔ̃]	BLOND	ON	[ŋ]	AGNEAU	Non Codé
SEMI-VOYELLES			* désignent les consonnes voisées		
[j]	YEUX	J -	[ɲ] est considéré comme décomposé en [n] + [j] (donc ~N + J-)		
[y]	HUILE	U -			
[w]	OUI	W -			

CHAPITRE II

CHAPITRE II

ANALYSE DU SIGNAL DE LA PAROLE EN VUE DE SA RECONNAISSANCE

II-1 INTRODUCTION

L'analyse acoustique est une partie importante dans le traitement que subit le signal de la parole pour pouvoir réaliser un système de reconnaissance de la parole.

Le but de cette analyse est d'extraire les paramètres pertinents (les énergies dans les bandes de fréquence, les coefficients cepstraux) qui caractérisent au mieux ce signal qu'est la parole.

Il existe plusieurs types d'analyse dont les plus importantes sont

- analyse spectrale par la transformée de fourrier rapide (F.F.T) ou par vocodeurs
- Analyse par la prédiction linéaire
- analyse cepstrale

Enfin citons les méthodes temporelles tels que le nombre de passage du signal par zero, la mesure de l'énergie et la fonction d'auto-correlation à court terme.

L'étude de ces différentes méthodes permet de choisir celle qui donne les meilleures performances le tableau ci-dessus indique le nombre de coefficients nécessaires pour une bonne restitution du signal de la parole pour chaque type de représentation.

.../...

types de coefficients	energie dans les bandes de frequence	coeff de la L.F.C	coeff du cepstre
nombre de coefficients	≈ 20	≈ 12	≈ 8
type de distance	Euclidienne	statistique	euclidienne

II-2 ANALYSE SPECTRALE

II-21 INTRODUCTION

L'analyse spectrale d'un signal échantillonné (numérique) consiste à le décomposer en une série infinie de sinusoides de fréquence f et d'amplitude $X(f)$. Cette décomposition est obtenue par la transformée de fourrier discrète (T.F.D). le calcul de la T.F.D est facilitée par les algorithmes de la F.F.T qui réduisent considérablement le temps de calcul.

L'analyse du spectre à court terme peut être obtenue à l'aide d'un vocodeur à canaux, enfin, les vocodeurs à formant effectue eux aussi l'analyse spectrale et détecte la fréquence des formants.

II 22- ANALYSE DU SPECTRE A COURT TERME PAR LA F F T

transformée de fourrier discrète (T.F.D)

soit une suite de valeurs discrètes d'un signal $x(k)$

la T.F.D est une suite de valeurs définies par $X(k)$

$$X(n) = \sum_{k=0}^{N-1} x(k) e^{j \frac{2\pi n}{N} k} \quad k=0, 1, \dots, N-1$$

la T.F.D inverse existe et a pour expression

$$x(k) = \frac{1}{N} \sum_{n=0}^{N-1} X(n) e^{-j \frac{2\pi k}{N} n} \quad k=0, 1, \dots, N-1$$

on peut l'écrire sous forme matricielle

on pose $W_N = e^{j\frac{2\pi}{N}}$ et $W_N^{rk} = (W_N)^{rk}$

$$\begin{bmatrix} X(0) \\ X(1) \\ \vdots \\ X(N-1) \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & W_N^1 & \dots & W_N^{N-1} \\ 1 & W_N^2 & \dots & \\ \vdots & \vdots & \ddots & \vdots \\ 1 & & & W_N^{(N-1)(N-1)} \end{bmatrix} \begin{bmatrix} X(0) \\ X(1) \\ \vdots \\ X(N-1) \end{bmatrix}$$

transformée de fourrier rapide (F.F.T)

si l'on considère la formule qui définit la T.F.D d'une suite de valeurs discrètes, il apparait qu'il convient de faire n^2 opérations (addition et multiplication) pour la calculer.

la F.F.T est une famille d'algorithmes qui permettent de réduire considérablement le nombre d'opération. En effet, l'algorithme de COOLEY-TUKHY nécessite $N \log_2 N$ opérations, soit un gain en temps-calcul de $N / \log_2 N$ ceci est évidemment intéressant lorsque N est très grand.

II-23 ANALYSE SPECTRALES PAR VOCODEUR.

1/ INTRODUCTION.

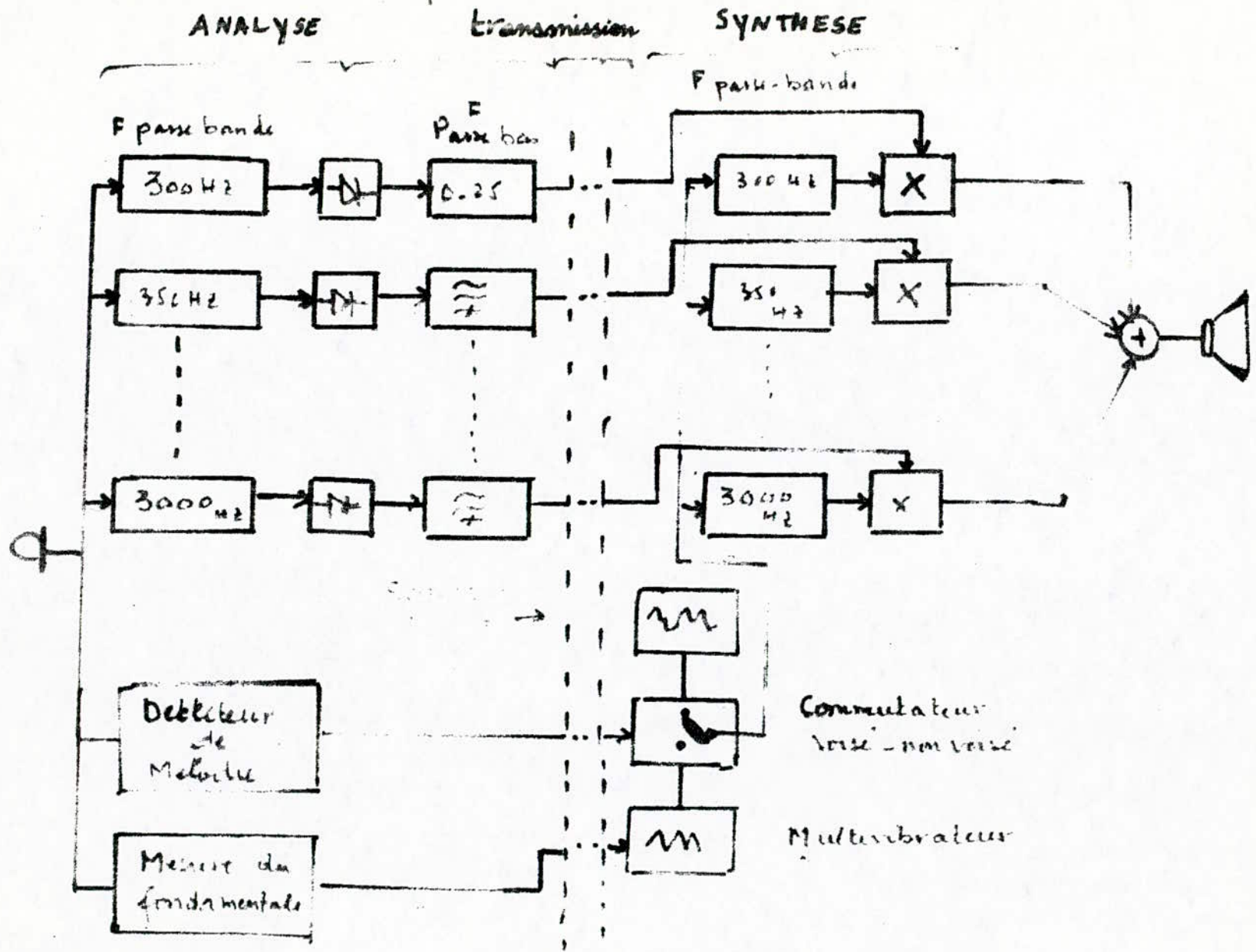
L'analyse spectrale de la parole peut être obtenue par les vocodeurs à canaux ou à formants. Ces derniers ont une fonction commune, la détection et la mesure de la fréquence du fondamentale. Chaque vocodeur se comporte d'un analyseur et d'un synthétiseur placés de part et d'autre du canal de transmission.

2-VOCODEUR A CANAUX.

En anglais le mot vocoder (voix coder) signifie codage de la parole l'analyseur du vocodeur à canaux assure deux fonctions distinctes.

-L'analyse spectrale; elle est obtenue par une batterie de filtres passe-bas recouvrant tout le spectre de la parole (300 à 3000 Hz). au nombre de 12 à 32 canaux. Chaque canal est constitué par un filtre passe bande suivi d'un redresseur et d'un filtre passe-bas permettant le lissage de l'énergie (fig 3).

.../...



fig(3)

VOCODEUR A CANAUX

-un détecteur de mélodie qui élabore des signaux à la fréquence du fondamentale. Ces signaux nous renseignent sur le caractère voisé ou non voisé du son analysé.

Le banc de filtres effectue donc un échantillonnage du spectre de la parole. En effet à la sortie de chaque filtre on peut mesurer la valeur de l'énergie dans chaque bande de fréquence toutes les 20 ms.

Ces énergies constituent des paramètres que nous pouvons utiliser dans les systèmes de reconnaissance automatique de la parole en général.

Le synthétiseur permet de reconstituer le message parlé selon un processus inverse de celui utilisé dans l'analyse, la parole synthétisée est intelligible mais le naturel de la voix est dégradé.

3-III VOCODEUR A. FORMANTS

L'analyseur du vocodeur à formant effectue une analyse spectrale du signal identique à celle du vocodeur à canaux et détecte les formants (fig 4) le synthétiseur par contre est constitué d'un ensemble de filtres résonnants dont la courbe de réponse globale en fréquence reproduit celle de conduit vocal, une source de bruit et une source d'impulsions périodique attaque ce circuit.

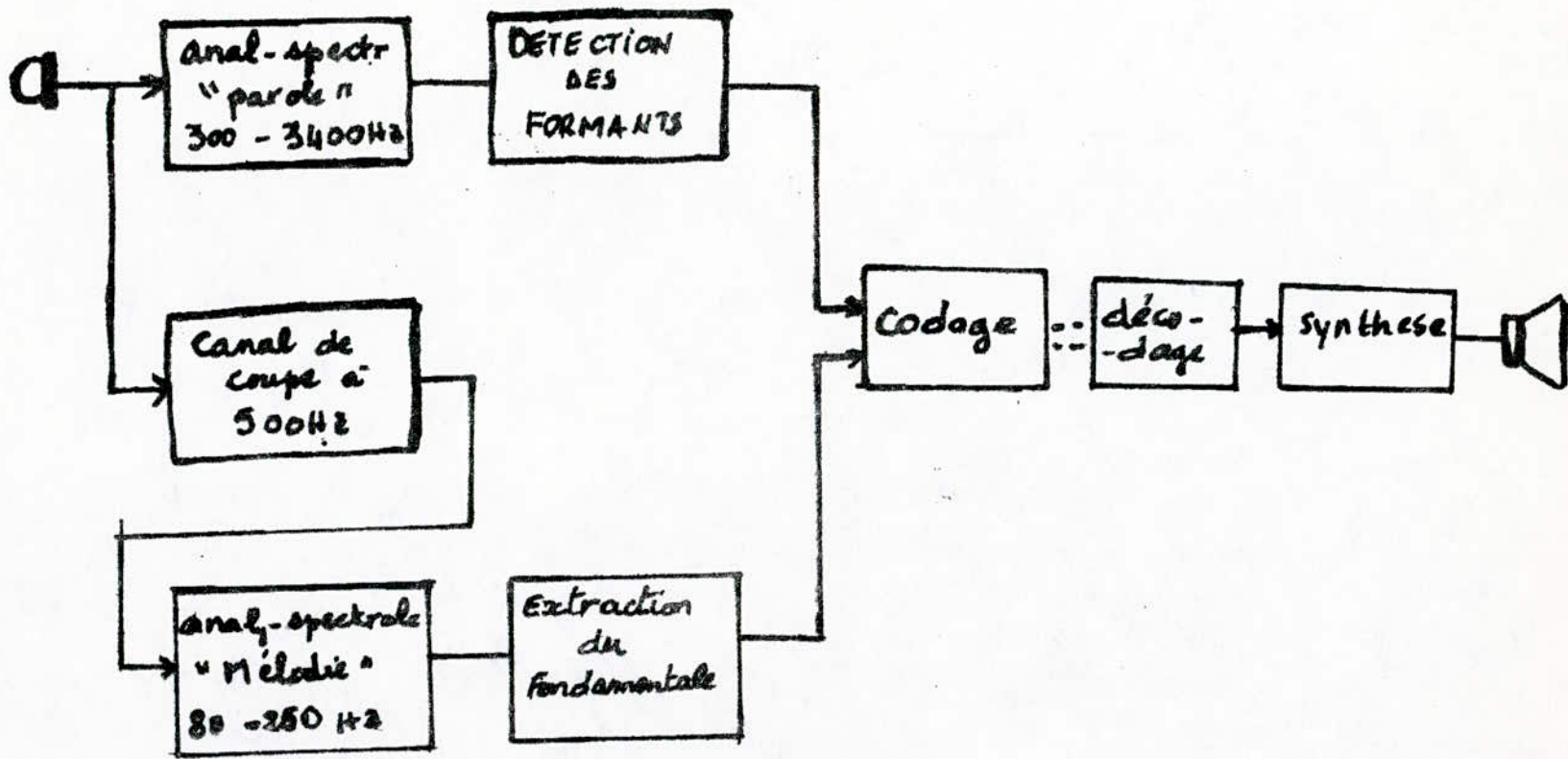
On remarque une certaine analogie entre le mode de production de la parole humaine et celle produite par ce circuit, ce qui laisse prévoir une meilleure synthèse de la parole que celle produite par le vocodeur à canaux. En effet, ce dernier ne tient pas compte des caractéristiques du conduit vocal, ni des contraintes qu'il impose. Malheureusement l'analyse par vocodeur à formant reste difficile à effectuer pour deux raisons :

- L'apparition de nombreux parasites sur le spectre provoque des erreurs importantes
- Le premier maximum peut être confondu avec le pitch.

4 CONCLUSION

Si les vocodeurs à formants donnent de bons résultats, il reste que la détection automatique des formants et en temps réel est difficile, c'est la raison pour laquelle les vocodeurs à formant ne sont pas utilisés fréquemment dans la synthèse de la parole.

.../...



Fig(4)

VOCODEUR A FORMANTS

II-3- ANALYSE PAR LA PREDICTION LINEAIRE (L.P.C)

II-31- INTRODUCTION

La prédiction linéaire est une technique d'estimation des paramètres de base de la parole (pitch, formants...) dont le critère d'optimisation est la minimisation de l'erreur quadratique moyenne.

L'analyse par la prédiction linéaire consiste à trouver un modèle à un signal original, ce modèle doit se rapprocher le plus de ce signal.

Chaque échantillon est estimé par une fonction linéaire des échantillons qui le précède immédiatement. Ce qui se traduit par la relation suivante:

$$S(n) = \sum_{k=1}^p a_k S(n-k) \quad (1) \quad a_k : \text{coefficients de la L.P.C}$$

Les avantages de la prédiction linéaire sont:

- une représentation du signal de la parole dans un espace de dimension réduite (de l'ordre de 10).
- une estimation précise des paramètres de base de la parole.
- une vitesse relative de calcul de ces paramètres.

II-32 - PRINCIPE DE L'ANALYSE PAR LA L.P.C

Pour mieux comprendre le principe de l'analyse par la L.P.C, on fait appel au modèle de la production de la parole par le système de la figure (5).

Ce système est excité par une source d'impulsions ou de bruit blanc.

D'une manière générale on peut exprimer tout signal temporel en terme de modèle prédit et d'un signal d'erreur

$$\text{on a } S(n) = \sum_{k=1}^p a_k S(n-k) + \sigma e(n) = \hat{S}(n) + \sigma e(n) \quad (2)$$

où σ est une constante d'adaptation d'énergie

la transformée en z de (-) est:

$$S(z) = \hat{S}(z) \left[\sum_{k=1}^p a_k z^{-k} \right] + \sigma e(z) \quad (3)$$

le modèle du filtrage est le suivant:

$$S(z) = \frac{\sigma}{1 + \sum_{k=1}^p a_k z^{-k}} E(z) = H(z) \cdot E(z) \quad (4)$$

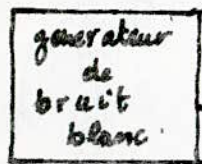
parole de Melodie



Paramètres du filtre numérique

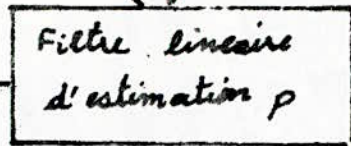


Sonores



sonde

Amplitude



$$S_m = \sum_{k=1}^p a_k S_{m-k}$$



$U(n)$



S_n

échantillon de parole

MODELE NUMERIQUE DE LA PRODUCTION DE LA PAROLE

Fig (5)

donc le modèle de la production de la parole, $H(Z)$ représente la fonction de transfert du conduit vocal.

II.33 METHODES DE CALCUL DES COEFFICIENTS DE LA L.P.C

Dans la théorie il existe une quantité illimitée de façons à calculer les coefficients de la LPC dépendant du critère d'approximation désiré entre le modèle et le signal.

dans notre cas on veut minimiser la quantité définie par

$$E = \sum_n (S(n) - \hat{S}(n))^2 = \sum_n \left(S(n) - \sum_{k=1}^p a_k S(n-k) \right)^2 \quad (5)$$

cette quantité est minimale si

$$\frac{\partial E}{\partial a_k} = 0 \iff \sum_{k=1}^p a_k \sum_n S(n-k) S(n-i) = \sum_n S(n-i) S(n) \quad (6)$$

le système linéaire d'équation se résoud soit par la méthode de covariance, soit par l'auto-corrélation. Le choix de la méthode dépend de l'intervalle d'analyse du signal.

II-32-1 METHODE DE COVARIANCE.

Cette méthode fait les suppositions suivantes

- le signal est de fini par $P+N$ échantillons où p est l'ordre du prédicteur et N la taille de l'échantillon estimé.
- Un échantillon est prédit par les p échantillons qui le précède. Ceci est valable exclusivement pour les N échantillons successifs.
- l'erreur est minimisée exclusivement sur les N échantillons.

Le système (6) se traduit par

$$\sum_{k=1}^p a_k \phi_{ik} = \phi_{i0} \quad i=1,2,\dots,p$$

$$\text{avec } \phi_{ik} = \begin{matrix} S(n-i) \cdot S(n-k) & i=1,\dots,p \\ & k=1,\dots,p \end{matrix}$$

la forme matricielle est:

$$\begin{pmatrix} \phi_{11} & \phi_{12} & \dots & \phi_{1p} \\ \phi_{21} & \phi_{22} & \dots & \phi_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{p1} & \phi_{p2} & \dots & \phi_{pp} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{pmatrix} = \begin{pmatrix} \phi_{10} \\ \phi_{20} \\ \vdots \\ \phi_{p0} \end{pmatrix}$$

La matrice de gauche est la matrice de covariance du signal $S(n)$, la solution des coefficients de la L.P.C est réalisé par la méthode de CHOLESKY, cet algorithme nécessite $\frac{p^2}{2}$ mots mémoires et effectue $p^3/6$ produits.

II_322 METHODE D'AUTO-CORRELLATION

Cette méthode fait les suppositions suivantes:

- le signal est nul à l'extérieur de l'intervalle considéré. Ceci est réalisé en multipliant le signal par une fenêtre temporelle de largeur N .
- chaque échantillon est prédit par ses p échantillons qui le précèdent et ceci pour tout le temps. ($n \in]-\infty, +\infty[$)
- l'erreur quadratique totale entre le signal fenêtré et le modèle est minimisé de $-\infty$ à $+\infty$

ces considerations donnent les équations normales d'auto-corrélation à partir de (6).

$$\sum_{k=1}^p a_k R_{|i-k|} = R_i \quad i = 1, 2, \dots, p$$

$$R_i = \sum_{n=0}^{N-1-i} S(n) S(n+i)$$

Les équations d'auto-corrélation ont la forme matricielle suivante

$$\begin{pmatrix} R_0 & R_1 & \dots & R_{p-1} \\ R_1 & R_0 & & R_{p-2} \\ \vdots & \vdots & \ddots & \vdots \\ R_{p-1} & R_{p-2} & \dots & R_0 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{pmatrix} = \begin{pmatrix} R_1 \\ R_2 \\ \vdots \\ R_p \end{pmatrix}$$

la matrice de gauche est carrée , en plus d'être symétrique, a les mêmes valeurs dans les lignes parallèles à la diagonale principale.

Une matrice de ce type est une forme particulière de matrice de TOEPLITZ. La solution des équations est donnée par l'algorithme de Lévinon. Ce dernier nécessite $2P$ mots mémoires et effectue $P(p+1)$ produits

CONCLUSION

En synthèse l'analyse par la L.P.C donne de bons résultats pour le calcul des paramètres de base de la parole.

En reconnaissance ces résultats sont discutés à cause des hypothèses prises au départ sur la stationnarité du signal.

La méthode de covariance assume que le signal de la parole est encore non stationnaire pour l'intervalle considéré dans le calcul des coefficients a_k . Cette méthode n'assure pas la stabilité du modèle.

La méthode d'auto-corrélation simplifie ces difficultés, en effet elle considère le signal stationnaire, le fenêtrage du signal avant l'estimation assure la stabilité du modèle. Son avantage c'est d'utiliser un minimum d'information pour le calcul des coefficients de la L.P.C

41/ INTRODUCTION

Les systèmes linéaires occupent une place importante dans le traitement des signaux car ils sont décrits par des relations mathématiques simples et d'utilisation souples, donc ils sont faciles à analyser. Malheureusement, ce n'est pas le cas pour les systèmes non linéaires, il n'est pas aisé et même difficile à représenter mathématiquement.

Le traitement homomorphique est le plus adapté pour l'analyse de ces signaux non linéaires, basé sur le principe de superposition généralisé.

42/ SUPERPOSITION GENERALISEE

on dit qu'un système est linéaire si et seulement si

$$L\{x_1(k) + x_2(k)\} = L\{x_1(k)\} + L\{x_2(k)\} \quad (1)$$

$$L\{a(x(k))\} = a L\{x(k)\} \quad (2)$$

le filtrage linéaire peut être utilisé pour séparer deux signaux occupants deux bandes de fréquence différentes, il ne peut être utilisé que si le signal recherché est combiné par addition.

Ex: $x(k) = x_1(k) + x_2(k)$ et que le signal utile est $x_1(k)$. Il suffit d'utiliser un filtre passe-bande qui laisse passer uniquement la bande de fréquence de $x_1(k)$

La question qui se pose maintenant est la suivante:

comment extraire le signal $x(k)$ s'il était combiné avec un autre signal par multiplication ou par convolution?

$$\text{ex: } x(k) = x_1(k) \cdot x_2(k) \quad \text{ou} \quad x(k) = \sum_{l=-\infty}^{\infty} x_1(l) \cdot x_2(k-l)$$

pour généraliser les relations (1) et (2) on écrira:

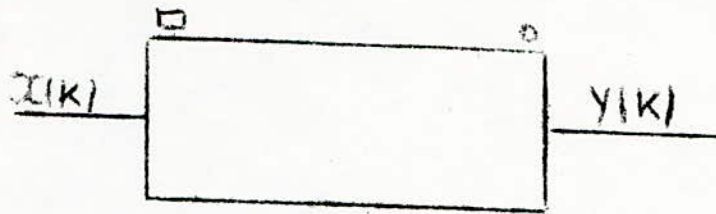
$$H(x(k)) \circ (x(k)) = H(x(k) \circ x(k))$$

$$H(c \circ x(k)) = H(x(k)) \circ c$$

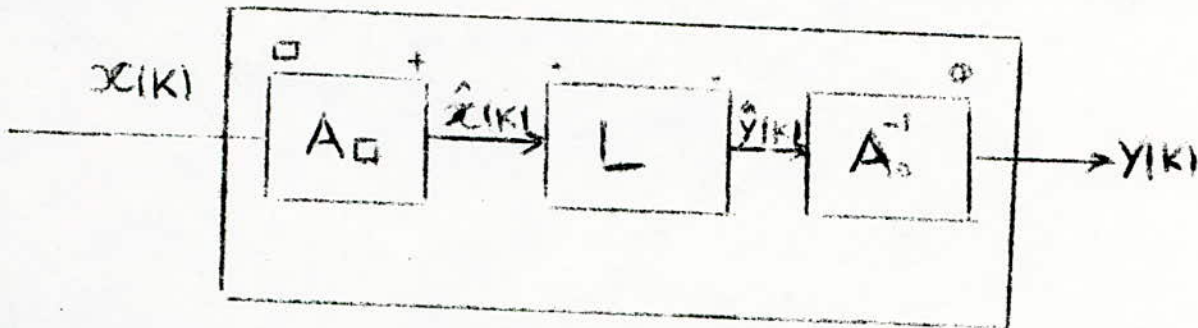
H représente l'opérateur du système

ainsi H est une transformée linéaire (structure algébrique d'un espace vectoriel). Ces systèmes sont représentés entre espaces vectoriels sont appelés homomorphiques.

un système homomorphique H avec la loi d'entrée et de sortie respectivement \square et \circ est représenté schématiquement par:



OPPENHEIM a démontré qu'il est possible de représenter un système homomorphe par 3 systèmes en série, dont celui du milieu est un système linéaire conventionnel, c'est la forme dite canonique des systèmes homomorphes (fig 1)



pour illustrer cette représentation, considérons les systèmes homomorphes multiplicatifs et convolutifs.

SYSTEME HOMOMORPHIQUE MULTIPLICATIF.

Les signaux d'entrée possible sont de la forme

$$x(k) = [x_1(k)]^{a_1} \times [x_2(k)]^{a_2}$$

a_1, a_2 sont des constantes arbitraires

le problème qui se pose est de trouver un espace où les signaux sont séparés, c'est à trouver un système qui satisfait la relation suivante

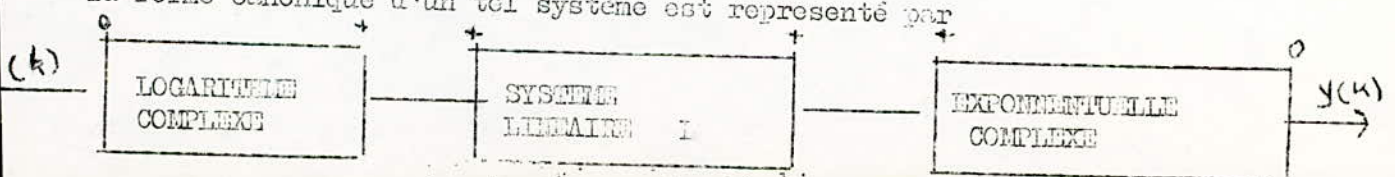
$$A [x_1^{a_1} \times x_2^{a_2}(k)] = a_1 A [x_1(k)] + a_2 A [x_2(k)]$$

$$\text{Log} [x_1^{a_1} \times x_2^{a_2}(k)] = a_1 \text{Log} [x_1(k)] + a_2 \text{Log} [x_2(k)]$$

dans le cas des signaux de la parole $x_1(k)$ & $x_2(k)$ peuvent être négatifs, on considère alors la fonction logarithme complexe

$$x(k) = |x(k)| \exp(j \arg x(k))$$

la forme canonique d'un tel système est représenté par



LES SYSTEMES HOMOMORPHIQUES CONVOLUTIFS

Les signaux sont tels que

$$y(n) = \sum_{k=-\infty}^{+\infty} h(n-k) \cdot x(k) = h(n) * x(n)$$

* représente le produit de convolution sur des signaux régulièrement échantillonnés. La déconvolution vérifie (D) l'équation:

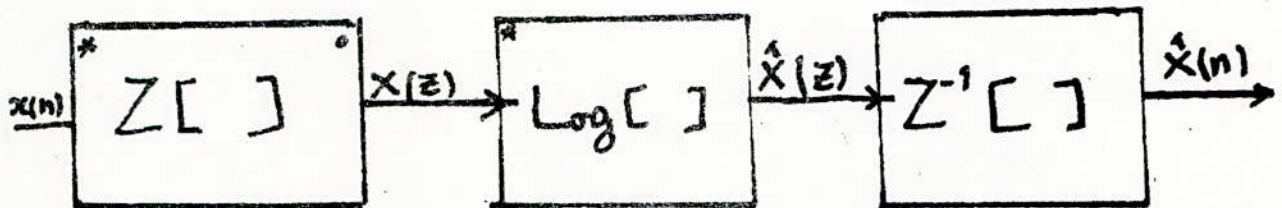
$$D(y(n)) = D(h(n) * x(n)) = D(h(n)) + D(x(n))$$

nous savons que le produit de convolution de deux signaux se transforme en une multiplication avec les transformées en z de chaque signal

$$Z(h(n) * x(n)) = H(z) \cdot X(z) = Y(z)$$

il suffit de prendre le logarithme complexe de Y(z) pour que le produit devienne une somme. Ainsi on peut représenter le système de déconvolution par la forme suivante:

$$\text{Log } Y(z) = \text{Log } h(z) + \text{Log } x(z)$$



déconvolution homomorphique

LE CEPSTRE

LE CEPSTRE EST LA TRANSFORMÉE DE FOURIER ^{inverse} du logarithme du spectre d'un signal.

Le cepstre $\hat{X}(k)$ d'un signal $X(k)$ est défini par :

$$\hat{X}(k) = -\frac{1}{j2\pi k} \int_{-1/2}^{1/2} \frac{1}{X(f)} \frac{d(X(f))}{df} \exp(j2\pi kf) df \quad k \neq 0$$

$$\hat{X}(0) = \int_{-1/2}^{1/2} \text{Log}|X(f)| df$$

PROPRIÉTÉ DU CEPSTRE

nous allons voir une propriété du cepstre que nous démontrerons pas soit une fonction de transfert ayant la forme suivante :

$$X(f) = A \cdot \frac{\prod_{n=1}^{Zc} (1 - d_n e^{-j2\pi f})}{\prod_{n=1}^{Pc} (1 - \gamma_n e^{-j2\pi f})} \cdot \frac{\prod_{n=1}^{Zc} (1 - \beta_n e^{j2\pi f})}{\prod_{n=1}^{Pc} (1 - s_n e^{j2\pi f})}$$

où $A > 0$ et $|d_n|, |\beta_n|, |\gamma_n|$ et $|s_n|$ est inférieur à 1
les coefficients cepstreux sont donnés par

$$\hat{X}(k) = \begin{cases} - \sum_{n=1}^{Zc} \frac{d_n^k}{k} + \sum_{n=1}^{Zc} \frac{\gamma_n^k}{k} & k > 0 \\ \sum_{n=1}^{Pc} \frac{\beta_n^k}{k} + \sum_{n=1}^{Pc} \frac{s_n^k}{k} & k < 0 \end{cases}$$

$$\hat{X}(0) = \text{Log}(A)$$

si le signal $x(k)$ est causal c'est à dire que:

$x(k)=0$ pour $k < 0 \Rightarrow \hat{x}(k)=0$ pour $k < 0$
 ceci se traduit par le fait que la fonction de transfert n'a ni zéro, ni pôles en dehors des cercle unite

$$X(f) = A \cdot \frac{\prod_{n=1}^L (1 - \alpha_n e^{-j2\pi f})}{\prod_{n=1}^R (1 - \delta_n e^{-j2\pi f})}$$

ECHELLE MEL DES COEFFICIENTS CEPSTRAUX.

On remarque les coefficients cepstraux décroissent en $\frac{1}{k}$, il suffit donc d'un nombre réduit pour caractériser le signal

l'échelle spectrale dans le calcul du cepstre est linéaire les études physiologiques et perceptives de l'oreille semblent indiquer que cette dernière est sensible à une échelle logarithmique ~~que~~ de la fréquence. c'est pour cette raison qu'on utilise une échelle quasi-logarithmique (échelle MEL) linéaire sur le premier KHZ et logarithmique au delà de 1 KHZ.

Les coefficients cepstraux sont obtenus à l'aide d'un vocodeur à canaux composés de 25 filtres passe bande de forme triangulaires (fig 7) Les 10 premiers filtres ont une largeur uniforme et les autres ont une largeur qui suivant une progression géométrique. Les coefficients cepstraux dans l'échelle MEL se calcule à partir de l'équation suivantes:

$$MFCC(n) = \frac{1}{NF} \sum_{k=1}^{NF} \text{Log } E(k) \cos(n(k-i)) \frac{\pi}{NF}$$

où NF: nombre de filtres triangulaire (en générale 25)

M le nombre de coefficients cepstraux

E (k) sont les reponses d'énergie des filtres triangulaire (voir fig 8)

CONCLUSION

Le cepstre a ces avantages suivants

- le spectre logarithmique rend visible les faibles densités spectrales dans les bandes de fréquence sensible à la perception
- la large gamme dynamique de la densité spectrale est réduite par la conversion logarithmique quantitativement. La conséquence de cette opération se représente dans le domaine temporel par une concentration de l'énergie autour de l'origine c'est à dire la plupart de l'énergie du cepstre se trouve dans la première dizaine de coefficients.

.../...

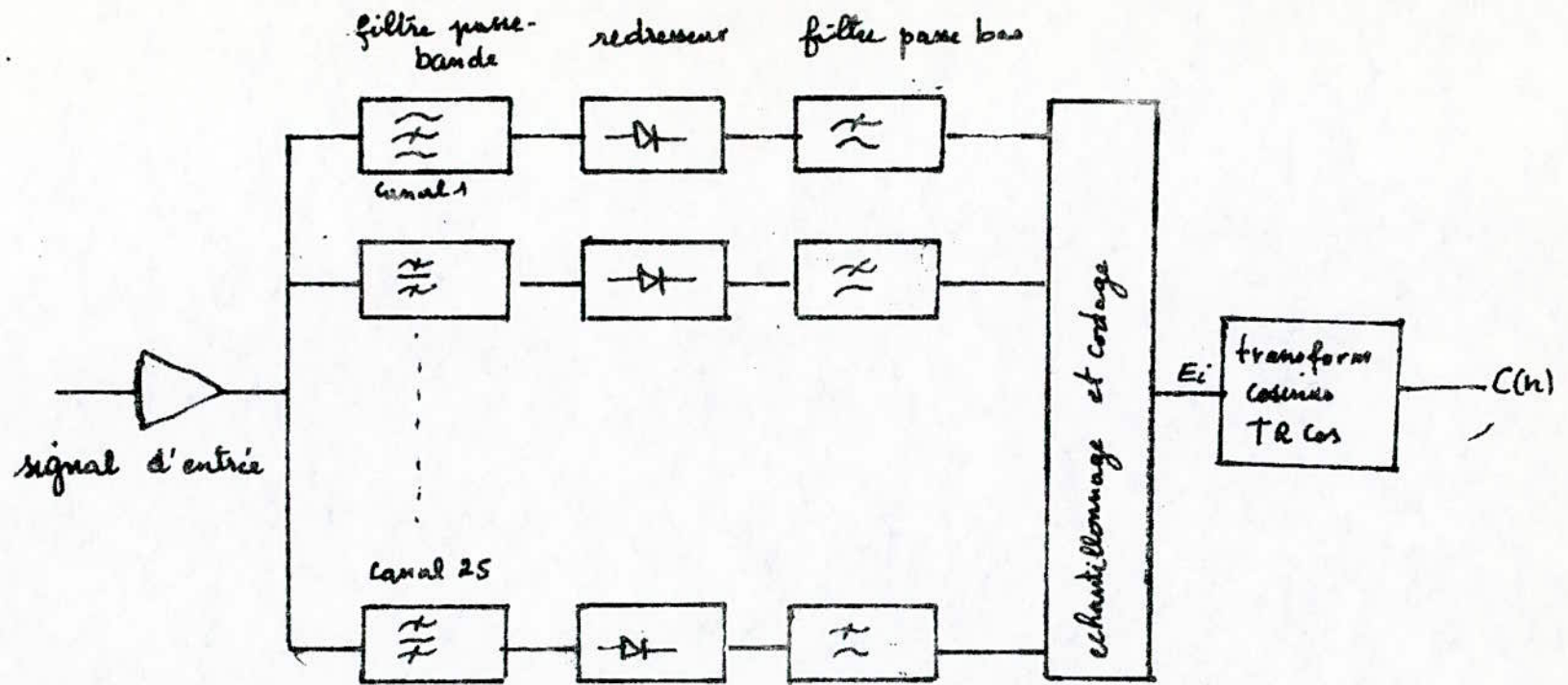
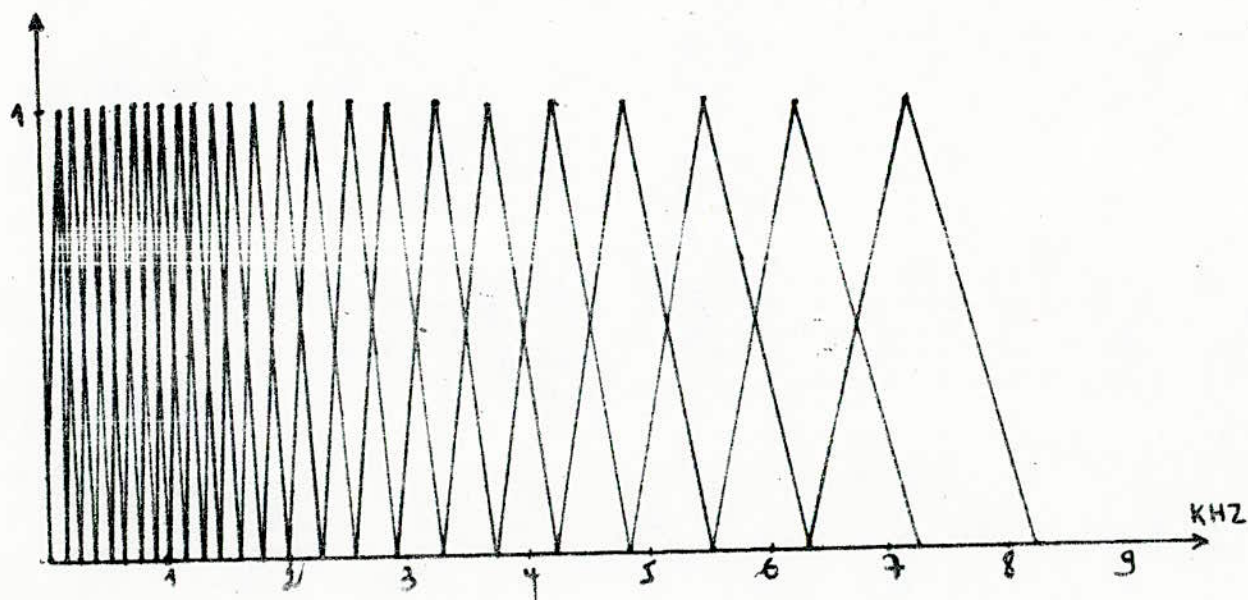


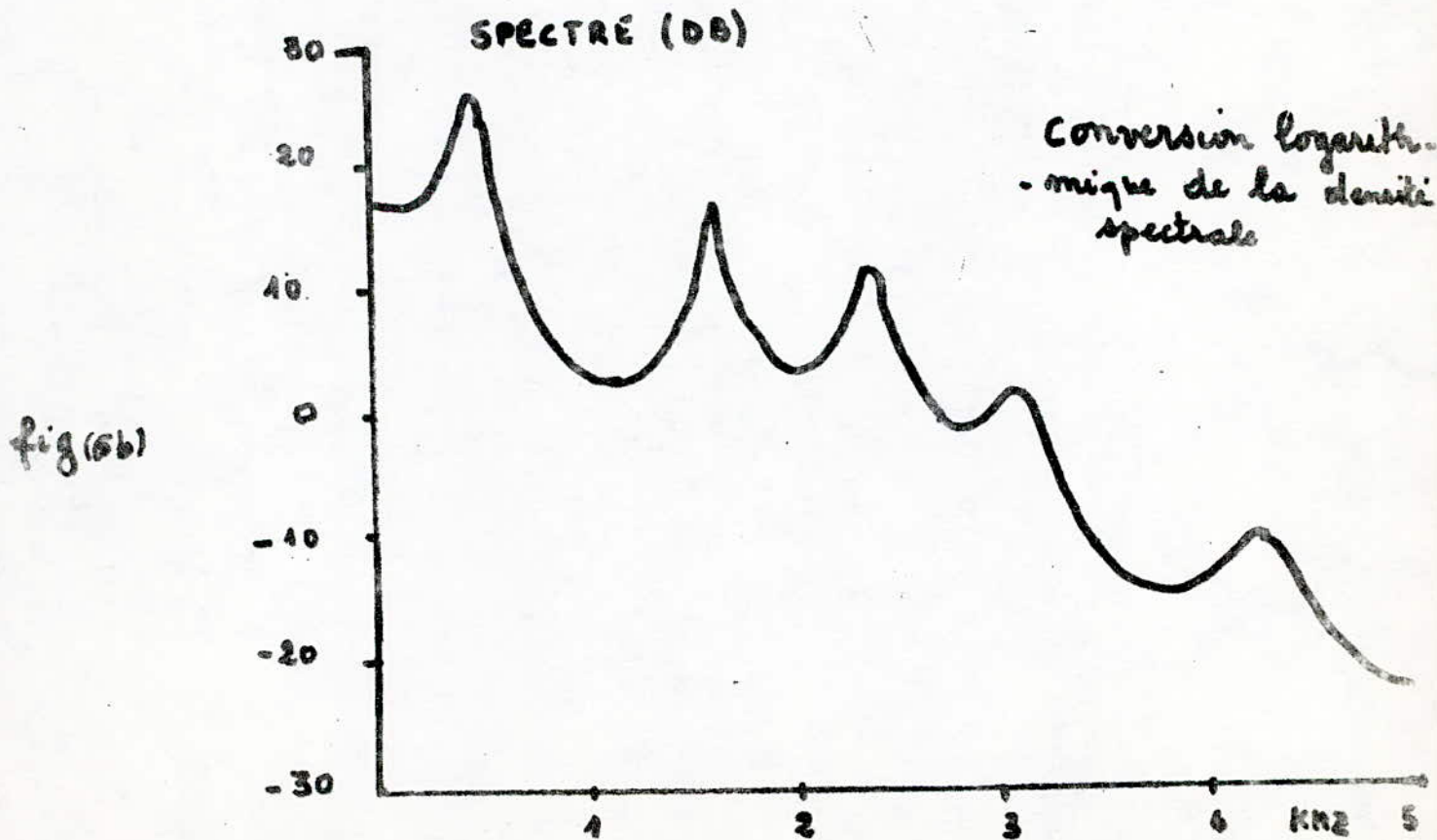
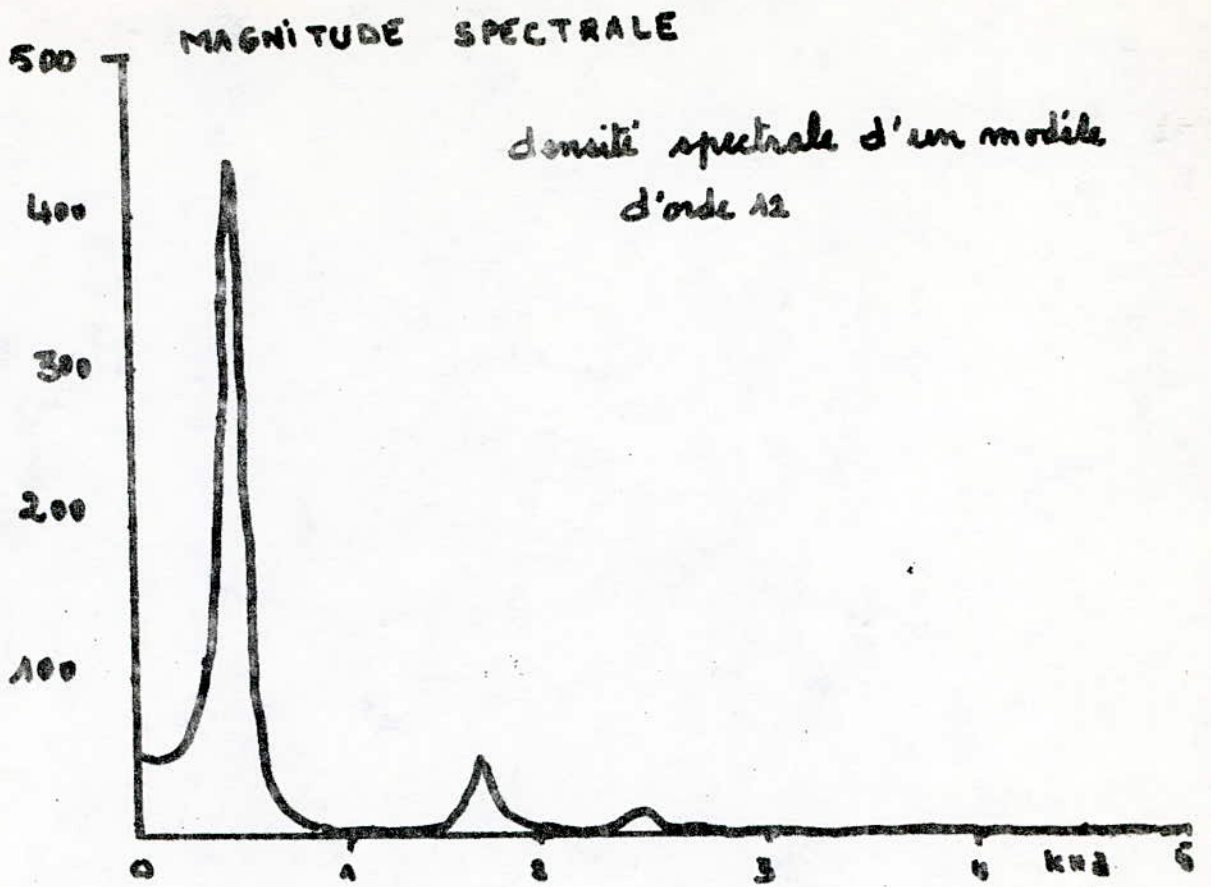
fig (7)

OBTENTION DES COEFFICIENTS CEPSTRAUX
MEL PAR UN VOCODEUR A CANAUX



Fig(8) Bancs de FilTres
Eriongulaires

c) la conversion logarithmique souleve les formants de tous faibles énergie localisés fréquemment dans la bande fréquentielle perceptible (voir fig 6) le cepstre a néanmoins un inconvenient l'utilisation de la distance cepstrale comme indice de similitude dans les systèmes de reconnaissance, nécessite un temps de calcul considerable.



II-5 ANALYSE TEMPORELLE

Ces méthodes d'analyse privilégient l'évolution temporelle de certains paramètres dont les principaux sont la fréquence du fondamentale f_0 , les formants F_1, F_2, \dots , et certains coefficients de la prédiction linéaire.

II-51 LA MESURE DE L'ÉNERGIE

L'énergie E d'un signal $x(n)$ échantillonné est donnée par la relation suivante

$$E = \sum_{n=-\infty}^{+\infty} x^2(n)$$

on considère généralement des intervalles de temps très courts (environ 20 ms) pour lequel le signal de la parole est stationnaire l'énergie dans cet intervalle est :

$$E = \sum [h(m)x(n+m)]^2$$

où $h(m)$ est une fenêtre temporelle dont le but est de donner le moins d'importance aux échantillons éloignés dans le temps cette énergie E est un paramètre simple et utile pour détecter les silences ou séparer les syllabes.

II-52 LE NOMBRE DE PASSAGES PAR ZÉRO DU SIGNAL

La méthode pour calculer le nombre de passage par zéro est facile à mettre en oeuvre, en effet, pour un signal discret, il suffit de vérifier la relation suivante : $x(n) \cdot x(n-1) < 0$

ce paramètre est utilisé pour la mesure de la fréquence du fondamentale et pour estimer la fréquence des formants. Cette méthode présente certains avantages,

- elle est indépendante de l'amplitude du signal

- elle est très rapide et pourrait fonctionner en temps réel

- elle permet de séparer les fricatives non voisées des phonèmes voisés

II-53 LA FONCTION D'AUTO CORRÉLATION À COURT-TERME

Cette fonction est définie pour un signal $x(t)$ échantillonné par

$$\Phi(m) = \lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{n=-N}^{+N} x(n)x(n+m)$$

le signal de la parole étant considéré stationnaire sur un intervalle court (15 à 20 ms), on définit alors la fonction d'auto-corrélation à court terme par

$$\left\{ \begin{array}{l} \phi_N(m) = \frac{1}{N} \sum_{l=0}^{N-|m|-1} x(l) x(l-m) \\ \phi_N(m) = 0 \text{ pour } |m| \geq N \end{array} \right. \quad m = 0, \pm 1, \dots, \pm(N-1)$$

si le signal $x(t)$ est périodique alors $\phi_N(m)$ est aussi périodique ceci laisse entrevoir un moyen de détecter les périodicités qui peuvent être invisibles au premier examen. Cette fonction permet aussi de calculer les coefficients de la prédiction linéaire.

CHAPITRE III

CLASSIFICATION

1 INTRODUCTION

Nous avons ^{vu} au chapitre précédent les différents types d'analyse. Chaque méthode nous donne un certain nombre de paramètres, ces derniers serviront de données pour la classification afin de constituer un système de reconnaissance automatique de la parole.

La classification est le rangement d'entités en groupes prédéfinis. Donc l'espace dans lequel évolue ou sont contenus les éléments à classer est un espace dont les caractéristiques sont entièrement connues. La classification quant à elle est une famille d'algorithmes dont le but est, partant d'un espace totalement inconnu, de définir des classes à partir de l'ensemble des éléments qui le composent.

L'ensemble des données (les distances) concernant les éléments de cet espace inconnu suffit à le définir. La classification s'attache à recouvrir la structure profonde des données. ces données initiales sont complexes et en quantité importante, donc difficile à cerner.

La définition des classes par la classification permet dans une étape suivante et en supposant que ces classes recouvrent tout l'espace d'étude, de procéder à un simple classement, la première étape de la classification dite d'apprentissage permet de choisir des représentants pour chaque classe. Ces derniers constitueront le dictionnaire des références, ensuite vient la deuxième étape dite de reconnaissance qui permet d'identifier un élément inconnu c'est à dire retrouver dans le dictionnaire d'éléments qui lui ressemble le plus.

III-2 Les principes de la classification numérique.

Les principes de la taxonomie (étude de la classification, ses bases, ses opérations et ses règles) ont été définis lors de l'introduction des procédures numériques dans les méthodes de groupement d'entités en classes. Ces procédures basant leur décision sur des caractères codés ou numériques lors des observations expérimentales, obligent à formuler des hypothèses concernant la nature des entités groupées.

.../...

0 CLASSIFICATION

La classification comme nous l'avons définie précédemment est une procédure mais le résultat de cette opération est souvent appelé classification. à priori on peut classer les éléments d'un espace inconnu de plusieurs manières, par conséquent on doit avoir une idée claire quant au résultat qu'on veut obtenir, il s'agit là de sélectionner efficacement les propriétés utilisées comme base de notre classification.

TAXONOMIE NUMÉRIQUE.

L'ouvrage de Sneath P et Sokal R " THE principales of numerical taxonomy" a marqué l'introduction des méthodes de classification automatique. toutes ces méthodes démarrent sur le choix d'un indice de similitude entre les éléments à classer et d'un algorithme qui établit une classification ou une hiérarchie de classification susceptibles de respecter les relations de similitude la base fondamentale de ces méthodes est la conversion de l'information portée sur les objets en quantités numériques et ceci afin de pouvoir estimer le coefficient de similitude.

Les principes fondamentaux de la taxonomie numérique peuvent se résumer ainsi.

1. Une classification est d'autant meilleure que la quantité des caractères sur laquelle elle est fondée est grande et que l'information portée pour chacun des caractères considérés est aussi grande que possible.

2. Pour créer des classes naturelles, on doit affecter une même pondération pour chaque caractère considéré.

3. La classification est basée sur un indice global de similitude

4. La similitude globale entre deux éléments est une fonction des similitudes entre chacun des caractères simultanément comparés.

II. ETUDE DE LA DISTANCE

Nous avons vu que toutes les méthodes de classification automatique démarrent par le choix d'un indice de similitude, celui-ci guide les décisions de l'algorithme de classification l'indice de dissemblance est en général une distance calculée à partir d'une des représentations paramétriques.

En général la fonction distance $d(x, y)$ entre deux éléments d'un espace de dimension n est définie par l'application

$$\begin{array}{ccc} \mathbb{R}^n \times \mathbb{R}^n & \xrightarrow{d} & \mathbb{R}^+ \\ (x, y) & & d(x, y) \end{array}$$

Cette distance doit avoir les propriétés suivantes:

$$d(x,y)=0 \iff x=y$$

$$d(x,y) \geq 0$$

$$d(x,y) = d(y,x) \quad \text{symétrie}$$

$$d(x,z) \leq d(x,y) + d(y,z) \quad \text{Inégalité triangulaire}$$

Dans le cas de la parole, cette distance doit avoir une interprétation spectrale et doit être aussi simple à calculer.

$$\text{soient } x = (\alpha_1, \alpha_2, \dots, \alpha_n)$$

$$y = (\alpha'_1, \alpha'_2, \dots, \alpha'_n)$$

$$d(x,y) = \sum_{i=1}^n |\alpha_i - \alpha'_i| \quad \text{distance de Minkowski}$$

$$d(x,y) = \left[\sum (\alpha_i - \alpha'_i)^2 \right]^{1/2} \quad \text{" " " " " euclidienne}$$

$$d(x,y) = \text{Max } |\alpha_i - \alpha'_i| \quad \text{" " " " " Chebychev}$$

CE sont les principales distances

III 4 L' APPRENTISSAGE

4 I INTRODUCTION

Le but de l'apprentissage comme nous l'avons signalé au début du chapitre, est de constituer un dictionnaire d'éléments prototypes capable de recouvrir l'espace étudié.

A Partir de la séquence de vecteurs formes observés lors de l'analyse du processus, on détermine des classes à l'aide d'un algorithme de classification.

On distingue trois(3) types :

- L'apprentissage supervisé: il y'a une suite d'échantillons classés dont on connaît la répartition de tous les éléments (échantillons d'apprentissage classés).
- L'apprentissage non supervisé :les estimations des paramètres se forment d'après les échantillons non classés.
- S'il se produit durant l'apprentissage une modification des paramètres(ou de la structure de l'algorithme de classification) optimisant ces algorithmes au point de vue du critère de qualité choisi.Ces dispositifs sont dits adaptatifs .

4-2 Les méthodes d'apprentissage non supervisées

Dans ces cas là on ignore le nombre et la nature des classes , il s'agit donc de grouper les vecteurs formes en sous ensembles. Nous citons les principales méthodes

Méthode des centroïdes:

On sait par exemple qu'il existe 3 classes à déterminer, on choisit 3 vecteurs formes quelconques de la séquence. Ces derniers serviront de premières estimations des centroïdes de ces trois classes. On associe alors à chaque centroïde les éléments qui lui sont proches, on ainsi coupé notre espace en trois domaines. On détermine alors les centres de gravités de ces domaines, puis on classe à nouveau les formes,

ce qui donne trois nouveaux domaines.

On démontre qu'au bout d'un fini nombre fini de cycles, les trois domaines et leurs centres cessent d'évoluer. On a ainsi isolé les trois classes.

Méthode de groupement en chaîne

Elle ne suppose à priori aucune connaissance sur le nombre de classes. Lorsque les amas sont assez denses et séparés on peut prendre les formes une à une dans un ordre quelconque, la première étant prise comme spécimen de la première classe, on calcule la distance qui la sépare de la deuxième forme. Si cette distance excède un certain seuil choisi à l'avance, on commence une deuxième classe, dans le cas contraire la deuxième forme est ajustée à la première classe. On procède ainsi jusqu'à l'épuisement de toutes les formes.

Méthode utilisant une mesure de similitude fondée sur le nombre de voisins communs :

soient n points X_1, X_2, \dots, X_n

on dresse pour chacun des points la liste des k plus proches voisins chaque point est à lui même son plus proche voisin d'ordre zéro; Les k autres voisins sont rangés par ordre de distance croissante à X_i . On considère toutes les paires (X_i, X_j) , on regarde d'abord si X_i figure parmi les k plus proches voisins de X_j et réciproquement. Si cette double condition est vérifiée et si en outre X_i ou X_j ont en commun un nombre k_s de voisins ($k_s < k$). On considère que X_i ET X_j appartiennent à la même classe, le nombre minimale k_s de voisins communs correspond à un seuil de similitude. Les classes ainsi obtenues, pouvant être à leur tour regroupées en super classes jusqu'à ce que le nombre et la densité des groupements obtenus convient à l'application envisagée.

METHODE D'APPRENTISSAGE SUPERVISEE

On connaît la répartition de tous les éléments c'est à dire le nombre et la nature des classes. On doit simplement estimer certains paramètres tels que la moyenne et la matrice de corrélation de chaque classe.

soient $m+1$ classes

$$\begin{matrix} X_1^0, & X_2^0 & \dots & X_{n_0}^0 \\ X_1^1, & X_2^1 & \dots & X_{n_1}^1 \end{matrix}$$

$$X_1^m, X_2^m, \dots, X_{n_m}^m$$

à partir de ces échantillons classés, on estime les paramètres de la densité de probabilité (moyenne et matrice de corrélation)

La moyenne de chaque classe est définie par:

$$A_k = \frac{1}{n_k} \sum_{i=1}^{n_k} X_i^k$$

et la matrice de corrélation par:

$$(M_k) = \frac{1}{n_k} \sum_{i=1}^{n_k} (A_k - X_i^k) (A_k - X_i^k)^t$$

APPRENTISSAGE ADAPTATIFS

Les formes à reconnaître peuvent subir des modifications au cours du temps, ainsi la voix du locuteur sous l'effet de la fatigue change.

Un certain nombre d'algorithmes permettent une adaptation de la machine à de telles dérives des formes, par conséquent, il faut

reajuster les centres de gravités des classes.

Par exemple: soit H_0 l'hyperplan séparateur de 2 centres de gravités G_0^1 ET G_0^2 respectivement de 2 sous-ensembles E_0^1 et E_0^2 dont les formes sont connues, au moment de l'apprentissage comme appartenant aux classes C_1 et C_2 . Lorsque une fournée de formes est présentée elle se trouve divisée par H_0 en 2 sous ensembles E_n^1 ET E_n^2 dont la machine déterminera les nouveaux centres de gravités G_n^1 et G_n^2 , d'où un nouvel hyperplan H_n^0 qui sera le plan médiateur de ces derniers. Le processus est repris dès qu'une nouvelle fournée se présente, on voit qu'il permet de suivre l'évolution du phénomène. L'intervention de l'être humain est limité au tout premier stade.

III-5 LA RECONNAISSANCE.

Il s'agit dans l'étape de reconnaissance de ranger un élément inconnu dans une classe se trouvant déjà dans le dictionnaire. On procède de plusieurs façons suivant la structure de ce dernier .

-Chaque classe est représentée par un seul élément prototype, on calcule alors toutes les distances séparant l'élément à classer à tous les représentants des classes . Le candidat retenu est celui dont la distance est minimale: c'est la méthode du plus proche voisin.

-Le dictionnaire est formé de plusieurs références par mot (classe) . La phase de reconnaissance fournit un ensemble de candidats ordonnés par leur distance au mot à reconnaître.

Les techniques Knn (les " k nearest neighbours", les k plus proches voisins) prend l'ordre des mots non plus référence par référence, mais en prenant les 2,3 ou k premières références de chaque mot Knn.

Les principales variantes des Knn sont :

- a)-en partant de la première référence vers les références de distance croissante, dès que le $K^{i\text{eme}}$ représentant d'un mot est atteint, ce mot est retenu comme candidat dans le rang où il apparait. C'est une méthode suffisante mais elle ignore l'ordre d'apparition et les distances des références d'les mots précédents la $K^{i\text{eme}}$.
- b)-l'ordre des mots candidats est obtenu par ordonnancement de la somme des distances des K premiers représentants de chaque mot, ^{ainsi} est pris en compte par une contribution à la somme des distances.

CONCLUSION

Le choix de la méthode de classification dépend essentiellement de la nature des éléments à classer et du résultat qu'on veut obtenir, l'opération essentielle dans les problèmes de reconnaissance de forme est l'attribution d'une classe (appartenant à un ensemble fini de classes d'assignation) à un élément dit d'entrée ou inconnu.

Il y'a plusieurs méthodes d'assignation dépendant du type d'apprentissage.

III-6. ANALYSE EN COMPOSANTES PRINCIPALES

61 INTRODUCTION.

On a vu au chapitre II qu'il est possible de représenter le signal de la parole dans plusieurs espaces. La question qui se pose maintenant est de savoir combien de composantes faut il pour avoir une représentation suffisante pour une classification efficace.

Le résultat de la paramétrisation du signal de la parole est un tableau de dimension (N, P) car à chaque intervalle d'analyse (10 à 20 ms) on a un jeu de paramètres (énergies dans les bandes de fréquence, coefficients cepstraux...) et N représente la durée du segment analysé.

Souvent la quantité $N \times P$ de données est très grande, l'analyse en composante principale inspiré des méthodes d'analyse de données permet de réduire la dimension du tableau (N, P) . Ce dernier est considéré comme étant le représentant des coordonnées de N points dans un espace à p dimensions affectés chacun de la masse unité.

L'idée de base de l'analyse en composantes principale est de projeter ce nuage de points sur des axes privilégiés (axes factoriels). Dans cette opération de projection il ne s'agit pas de conserver la particularité de chacun des points mais d'extraire de ces points ^{des} caractéristiques moyennes .

62 PRINCIPE DE L'ANALYSE EN COMPOSANTES PRINCIPALES.

soient $X_i \quad i=1, 2 \dots \dots N$ avec $X_i \in \mathbb{R}^p$

X_i représente les données à classer, il y a donc $N \times p$ données

soit D_{kl} la distance entre deux points M_k et M_l de l'espace \mathbb{R}^p

on remplace D_{kl} par une projection d_{kl} de telle manière à avoir les quantités suivantes:

$$\frac{1}{N} \sum_{k=1}^N D_{kL}^2 \quad \text{et} \quad \frac{1}{N} \sum_{k=1}^N d_{kL}^2$$

les plus proches possibles c'est à dire que ^{la} ~~une~~ projection apportera le minimum de déformation du nuage des points considérés initialement.

63 DETERMINATION DES AXES FACTORIELS

On considère l'occurrence étudiée, c'est à dire le tableau

$$X = \left\{ X_i / X \in \mathbb{R}^p, i = 1, 2, \dots, N \right\}$$

Comme étant un tirage particulier réalisé sur une variable aléatoire de \mathbb{R}^p dont la distribution statistique est à déterminer, l'estimation statistique de la moyenne de la distribution cherchée et de sa covariance conduit aux valeurs suivantes.

$$\mu = \frac{1}{N} \sum \left\{ X_i / X_i \in \mathbb{R}^p \right\}$$

$$\Sigma = \frac{1}{N} \sum \left\{ (X_i - \mu)(X_i - \mu)^t \right\}$$

les axes principaux d'inerties du nuage sont aussi les axes factoriels. Ceux ci sont déterminés d'après la matrice de covariance

soit u le vecteur unitaire de \mathbb{R}^p , les points X_i peuvent être considérés comme des vecteurs de \mathbb{R}^p . Les projections sur la direction de u des points du nuage X sont les points de la droite portée par u et dont les coordonnées sont $y_i = \langle u, X_i \rangle$ (produit scalaire de u par X_i).

$$\text{En notation matricielle } y_i = u^t X_i$$

Le premier axe factoriel est l'axe sur lequel les projections des points du nuage ont une dispersion maximale. La variance de projection $V(u)$ est

$$V(u) = E \left((y_i - \bar{y})(y_i - \bar{y})^t \right)$$

où $\bar{y} = u^t \bar{X}$ \bar{X} est la moyenne des projections

$$V(u) = \frac{1}{N-1} \sum_{i=1}^N (u^t (X_i - \bar{X}))^t u^t (X_i - \bar{X})$$

$$V(u) = u^t T u$$

On remarque que T n'est autre que la matrice de covariance de l'échantillon X .

Le premier axe factoriel est le vecteur propre de T correspondant à la plus grande valeur propre il s'agit de déterminer le vecteur unitaire u minimisant $V(u)$, en utilisant le multiplicateur de Lagrange, il faut donc maximiser le critère J défini par :

$$J = u^t T u - \lambda (u^t u - 1)$$

derivons J par rapport à u (dérivation matricielle).

$$\frac{\partial J}{\partial u} = 2Tu - 2\lambda u = 0 \Leftrightarrow Tu = \lambda u$$

On voit bien que la valeur du critère $V(u)$ est égale à la valeur propre λ , le maximum de $V(u)$ est obtenu pour la plus grande valeur de λ de même les axes factoriels successifs du nuage sont les vecteurs propres de la matrice T de covariance de X rangés dans l'ordre des valeurs propres décroissantes

6-4 PROPRIÉTÉS DES AXES FACTORIELS

La valeur propre λ_k rapportée à la somme des valeurs propres représente la part d'inertie ou de variance totale du nuage portée sur l'axe factoriel défini par le vecteur u_k . Les axes factoriels ont de ce fait une propriété intéressante pour la réduction des données. En effet, on peut réduire la dimension de l'espace initial R^D sans pour autant déformer le nuage de points. une autre propriété des axes factoriels est la décorrélation réalisée sur les coordonnées des points.

65 CONCLUSION.

Dans le cas de la reconnaissance automatique de la parole, trois axes factoriels suffisent à décrire tout le nuage des points. Ceci dit l'application de cette méthode diminue sensiblement le taux de reconnaissance.

CHAPITRE

IV

(méthode des K-moyennes)

1- Introduction

L'étude et l'analyse du signal de la parole nous ont permis d'extraire des paramètres déjà cités au chapitre II. Ces paramètres caractérisent les éléments d'un espace. L'algorithme des K-moyennes, directement inspiré des méthodes des nuées dynamiques, a pour but de définir des classes à partir de cet espace.

Cette étape de la classification précède celle de la reconnaissance. Le but de la première étape dite d'apprentissage est de munir l'espace descriptif du processus d'une structure de représentation. Cette dernière appelée dictionnaire, est capable de recouvrir l'espace défini par une structure de classe. L'étape suivante, celle de la reconnaissance, consiste à assigner un élément quelconque à une des classes. Cette assignation est réalisée par une comparaison d'un élément inconnu à la totalité des représentants des classes possibles et une décision indicatrice de la classe la plus proche de l'élément.

2- Nuées dynamiques

21- Introduction

La méthode des nuées dynamiques est une méthode de classification automatique sans professeur due à DIDAY. Elle permet de déterminer sur un espace donné une structure de classes. Cette méthode est itérative, l'itération du processus se poursuit jusqu'à la stationnarité.

22- principe de la méthode

La méthode consiste à déterminer sur une population donnée une partition en un nombre k de classes, k fixé a priori. A chaque classe est associée un noyau.

On agrège la partition à partir des noyaux, puis à partir des partitions obtenues, trouver de nouveaux noyaux capable de générer une partition meilleur que la précédente. On itère le processus jusqu'à ce qu'il devienne stationnaire.

23- fondement de la méthode:

Soient I l'ensemble à classer

P_k l'ensemble des partitions de I en k classes disjointes.

L_k l'ensemble des noyaux, un élément de L_k est la réunion des k noyaux.

On se donne deux applications f et g telles que:

$f(L)=P$ est la fonction qui à un noyau associe la partition P , cette application réalise la reconnaissance.

$g(P)=L$ est la fonction associant à une partition P l'ensemble des noyaux L elle réalise une opération d'apprentissage

D'une partition P_i on passe à la partition P_{i+1} par:

$$P_{i+1} = (f \circ g)(P_i) = (f \circ g)(P_0)$$

Soit D la distance d'un élément de L à un noyau $\alpha_i: D(X, \alpha_i)$

dans notre cas, I est formé de vecteurs, et la distance:

$$D(X_n, \alpha_i) = \alpha_i X_n$$

La fonction f associe tout point au noyau dont il est le plus proche.

$$f(L)=P \text{ avec } P_i = \{X \in I, D(X, \alpha_i) \leq D(X, \alpha_j) \quad j \neq i\}$$

La fonction associe à une classe de I le vecteur optimal. En choisissant comme critère J :

$$J_i = \sum_{j \in P_i} \alpha_i^2 X_j$$

le noyau optimal est celui qui associe à la classe

P_i le vecteur (calculé par l'algorithme de DURBIN) sur le centre de gravité de la classe.

DIDAY montre que l'algorithme converge si la distance R par laquelle les noyaux sont choisis vérifie la propriété d'être carrée. Sans préciser plus cette propriété, signalons que le critère J que nous adoptons définit précisément une distance carrée et esquissons la démonstration de la convergence de l'algorithme.

Le critère J_i est positif ou nul. Nous associerons à la partition P_i le critère

$$J = \sum_{i=1}^k J_i$$

ce critère vaut donc
$$J = \sum_{i=1}^k \alpha_i^2 \sum_{j \in P_i} X_j = J(\alpha, P)$$

montrons que:

$$J(\alpha^n, P^n) \geq J(\alpha^{n+1}, P^n) \geq J(\alpha^{n+1}, P^{n+1})$$

où α^n et P^n sont les noyaux et les partitions obtenues à la Nième itération. La première inégalité est liée à l'optimalité des noyaux α^{n+1} calculés sur P^n . Pour chaque classe P_i^n , le noyau α_i^{n+1} calculé sur P^n rend le critère J_i minimum:

$$J_i(\alpha^n, P^n) \geq J_i(\alpha^{n+1}, P^n)$$

La seconde inégalité vient de ce que les points de I sont ou bien dans la même classe entre la partition $P^{(n+1)}$ et la partition $P^{(n)}$ (et alors leur contribution au critère reste le même), ou bien s'ils ont changé de classe, leur contribution $\alpha_j^{n+1} X$ est devenue $\alpha_j^{n+1} X$ qui adiminue (le point est classé avec la classe j qui minimise $\alpha_j^{n+1} X$), et le critère global J ne peut que diminuer.

La suite des critères $J = J(\alpha^n, P^n)$ est décroissante, bornée par zéro, elle converge donc. En outre l'ensemble P_k est fini car l'ensemble I est fini et le critère J^n ne peut prendre qu'un nombre fini de valeurs.

La suite J^n est donc stationnaire, à partir d'un rang N les centres ne se modifient plus.

24- Conclusion:

Cette méthode de classification se justifie par sa simplicité, la convergence y est réalisée. Cela veut dire que la stationnarité est obtenue à partir d'uncertain rang puisque l'ensemble à classer est fini.

Cet algorithme de classification requiert assez de temps pour réaliser la convergence, néanmoins il est très performant sur le signal de la parole.

III.3. ALGORITHME DES K-MOYENNE

1. SEQUENCE D'APPRENTISSAGE

1.1 règles phonologique et coarticulation

Dans la conception de tout système de reconnaissance analytique, il est indispensable de tenir compte des influences qu'exercent les phonèmes sur leur voisins. Cela est dû à l'effet de coarticulation. Par exemple, le sonagramme révèle que le K est à plus haute fréquence dans (ki) que dans (kou); bien que l'oreille ne fasse pas la différence.

On envisage donc d'enregistrer plusieurs locuteurs, chacun des phonèmes plusieurs fois une série de mots, ainsi les règles décrivant les variations de prononciation affectant les phonèmes dans tel ou tel contexte sont prises en compte.

Implicitement les diverses prononciations figurent dans le dictionnaire.

1.2 description de la forme et prétraitement

A chaque période d'analyse on a un jeu de paramètres (vecteur), donc pour un segment donné, on a un ensemble de vecteurs (matrice). Le prétraitement consiste à faire une segmentation bruit-parole, puis une segmentation en phonème dans le cas de la reconnaissance des phonèmes.

Le prétraitement effectue une normalisation en temps, cela veut dire que pour une matrice temps-fréquence, on se ramène à un vecteur forme car chaque phonème a une durée propre qui dépend de sa nature et du locuteur qui le prononce.

Supposons par exemple numérisées les sorties d'un vocodeur à 12 canaux, soit n le nombre de prélèvements correspondant à la durée du mot analysé et soit k_i le nombre de prélèvements correspondant aux divers phonèmes i qui composent ce mot. On obtient donc un ensemble de matrices au nombre i et de dimension (k_i, n) qui décrivent tout le mot.

La normalisation permet de transformer la matrice i en un vecteur, ainsi chaque vecteur décrit un phonème. C'est sur ces données que va opérer l'algorithme des k-moyennes.

1-3 La classification

- principe de la methode

La méthode consiste à calculer le centroïde de toute la séquence d'apprentissage. Ce centre est perturbé, c'est à dire remplacé par un couple d'éléments voisins. On associe alors à chacun de ces deux éléments une classe dont on déterminera le centroïde. Les deux représentants obtenus sont à nouveau perturbés pour obtenir quatre éléments d'ou quatre nouvelles classes et quatre centroïdes.

A chaque étape le nombre de représentants est multiplié par deux, ainsi si k est une puissance de deux, on obtient au bout de $\log(k)$ itérations un dictionnaire de taille k .

La langue française se compose de 36 phonèmes mais il suffirait d'une trentaine pour qu'elle soit intelligible

Dans notre cas on prendrait 32 qui est une puissance de deux.

- Mise en oeuvre de l'algorithme

Supposant que la séquence d'apprentissage soit formée de N vecteurs

$$x_1, x_2, \dots, x_N$$

Soit $y(k) = \{y_i(k) \quad i=1, \dots, k\}$ le dictionnaire à la $k^{\text{ième}}$ itération.

La première étape de l'itération consiste à partitionner la séquence d'apprentissage en k classes $S_i(k)$ définies par:

$$x^h \in S_i(k) \Leftrightarrow d(x^h, y_i(k)) \leq d(x^h, y_j(k)) \quad j=1, \dots, k$$

L'erreur de quantification est calculée pour chaque classe $S_i(k)$

$$\text{par: } D_i(k) = \frac{1}{N} \sum_{x \in S_i(k)} (d(x, y_i(k)))$$

L'espérance mathématique est remplacée par une sommation, chaque vecteur de la séquence

x^1, x^2, \dots, x^N se voyant attribuer la probabilité $1/N$

Dans une deuxième étape, on calcule le centroïde ou représentant de chaque classe c'est à dire le vecteur C_i qui minimise $D_i(k)$. on calculera C_i en annulant les dérivées partielles de :

$$D_i(k) = 1/N \sum_{x \in S_i(k)} (d(x, y_i(k)))$$

$$\frac{\partial D_i}{\partial y_1} = 0, \dots, \frac{\partial D_i}{\partial y_m} = 0 \quad y = \{y_1, y_2, \dots, y_m\}$$

ce qui nous donne $y_j = 1/N \sum x_{ij}$

Le nouveau dictionnaire est formé de k représentants ainsi définis.

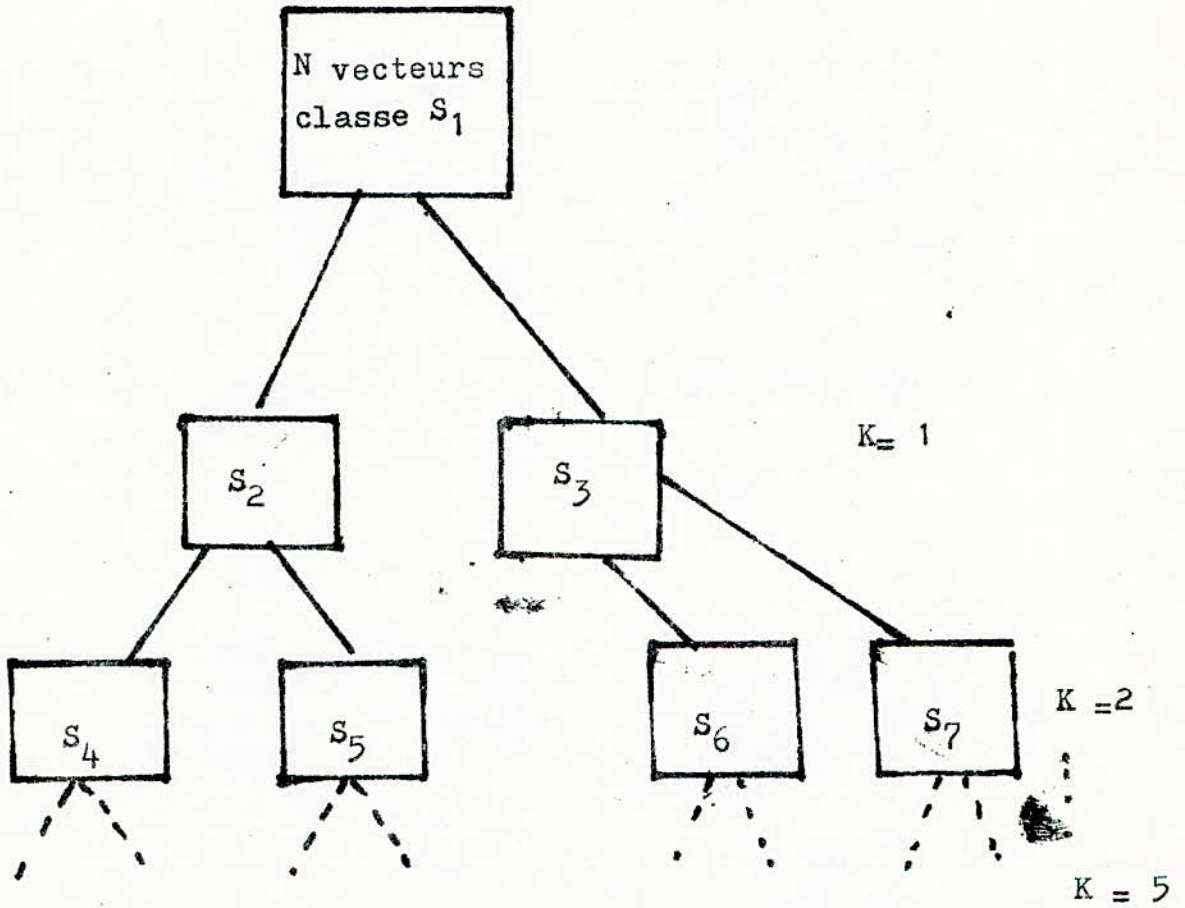
Le processus se poursuit jusqu'à ce qu'un état stationnaire soit atteint.

La somme des erreurs associées à chaque classe a atteint un minimum local. L'initialisation peut s'effectuer de différentes manières, la plus simple consiste à choisir aléatoirement k vecteurs de la séquence d'apprentissage.

Une fois que l'algorithme des k moyennes est élaboré, on sauvegarde par programmation toutes les classes structurées durant toutes les itérations.

Ainsi on confère au dictionnaire une structure en arbre comportant 63 classes (voir figure) •

SEQUENCE D'APPRENTISSAGE



STRUCTURE HIERARCHIQUE BINAIRE

ALGORITHME PRINCIPALE DES K-MOYENNES

LECTURE DU TABLEAU X

$S(1) = 1$; $S(2) = N+1$

Pour $I = 1$ à 33

 New $S(I) \leftarrow N + 1$

 FFaire

Appel centroïde ($1, N, C$)

$y(1) = C$

Pour $K = 1$ à 5

 Faire pour $I=1$ à 2 ($K+1$)

 Appel voisins ($I, y(200(K-1) + I-1), y_1, y_2$)

 Appel classes (I, y_1, y_2, S_1, S_2)

 Appel centroïde (S_1, S_2-1, C_1)

 Appel centroïde ($S_2, S(I+1) - 1, C_2$)

$y(2^k + 2I - 2) = C_1$

$y(2^k + 2I - 1) = C_2$

 New $S(2I - 1) = S_1$

 New $S(2I) = S_2$

 Fin faire

$S = \text{NEW } S$

FFaire

ALGORITHME DE CALCUL DU CENTROIDE DE LA
SEQUENCE D'APPRENTISSAGE

Centroide (ind1 , ind2 , C)

```
Pour j= 1      à m
  C1 ← 0
  Pour a = ind1 à ind2
    C1 = C1 + X(a,j)
  FFaire
  C(j) ← C1 / (ind2 - ind1 + 1 )
FFaire
```

COMMENTAIRES

ind1 , ind2 indices du début et de la fin de la séquence

C : centroide

ALGORITHME DE RECHERCHE DES PLUS PROCHES VOISINS DU CENTROIDE

Voisins(I, C, y₁, y₂)

y₁ ← X(S(I)) ; D_{min1} = d(X(S(I)), C)

y₂ ← X(S(I)+1) ; D_{min2} = d(X(S(I)+1), C)

Pour A = S(I) à S(I+1) - 1

Faire

D = d(X(A), C)

Si D < D_{min1}

alors D ← D
 D_{min1} ← D
 y₁ ← X(A)

Si D > D_{min1} et D < D_{min2}

alors D_{min2} ← D
 y₂ ← X(A)

Commentaires

S(I) : indice du début de la séquence

C : centroide de la séquence

y₁, y₂ les plus proches voisins de C

ALGORITHME DE CLASSEMENT

Classes(I, y₁, y₂, S₁, S₂)

Q ← 1, P = 0

Pour j = S(I) à S(I+1)-1

Faire

si d(X(j), y₁) < d(X(j), y₂)

/ X appartient à la classe S₁ /

alors

TAB(Q) ← X(j)

Q ← Q+1

sinon / X appartient à la classe S₂ /

TAB(S(I+1)-1-P) ← X(j)

P ← P+1

Fsi

S₁ ← S(I)

S₂ ← S(I)+Q-1

Q ← 1

Pour j = S(I) à S(I+1)-1

Faire X(j) = TAB(Q)

Q ← Q+1

FFaire

Commentaires

- y₁, y₂ voisins du centroïde
- I : indice du début de la séquence
- S₁ et S₂ classes associées à y₁ et y₂

III-4 RECONNAISSANCE DES PHONEMES

L'étape de reconnaissance consiste à identifier un phonème quelconque, ce phonème inconnu est comparé à la totalité des centroïdes (représentant de chaque phonème du dictionnaire).

La donnée de base de l'algorithme de reconnaissance est une distance. L'algorithme des k-moyennes ascendant au dictionnaire une structure hiérarchique binaire. Au lieu de calculer 32 distances et prendre l'élément le plus proche, on calcule simplement $2 \log_2(N)$ distances, d'où un gain en temps-calcul considérable.

L'algorithme de reconnaissance aura comme entrées un vecteur représentant le phonème à identifier et la distance D

PHONE (V , étiq , D)

Ind 1 = 2

Ind 2 = 3

Pour K = 1 à 5

Faire

$D_1 = \text{distance}(y(\text{ind } 1), V)$

$D_2 = \text{distance}(y(\text{ind } 2), V)$

si $D_1 < D_2$

alors $D \leftarrow D_1$

$\text{ind} \leftarrow \text{ind } 1$

$\text{ind}1 \leftarrow 2 \times \text{ind } 1$

$\text{ind}2 \leftarrow 2 \times \text{ind } 1 + 1$

Sinon

$D \leftarrow D_2$

$\text{ind}1 \leftarrow 2 \times \text{ind}1$

$\text{ind}2 \leftarrow 2 \times \text{ind}1 + 1$

Fin si

étiq = TRANS (ind - 31)

FIN

5- Conclusion

Pour diminuer le volume des calculs nécessaires à la recherche du meilleur représentant d'un vecteur, il est nécessaire de hiérarchiser les données en conférant au dictionnaire une structure en arbre.

L'algorithme que nous avons élaboré permet, à partir de la racine l'arbre, d'emprunter à chaque noeud la branche qui minimise la distortion.

LE PROGRAMME D'APPRENTISSAGE:

Le programme est écrit en P.L.1, nous avons choisi ce langage car il est puissant et universelle. Il est le plus adapté à notre problème, de plus il permet la sauvegarde de la structure en arbre du tableau des centroïdes de toutes les classes.

Les données de ce programme sont des vecteurs de dimension quelconque. Nous avons choisi arbitrairement cette dimension, mais elle peut être ajustée facilement en fonction du type de représentation obtenue lors de l'analyse des phonèmes.

Le type de distance utilisée est l'erreur quadratique pondérée. La pondération consiste à donner plus d'importance à certaines composantes du vecteur lorsque ces derniers portent plus d'informations que les autres.

PROGRAMME D'APPRENTISSAGE (NUEE)

```

1      0      NUUES : PROC ;
2      1      DCL SYSIN FILE STEAM INPUT ;
3      1      dcl SYSSPRINT FILE SREAM OUTPUT ;
4      1      DCL X(N,4) BIN FIXED(31) CTL ,
5      1      C(4) DEC FIXED (31,5) ,
6      1      Y(63,4) BIN FIXED(31) ,
7      1      (Y1(31),Y2(31)) BIN FIXED (31) ,
8      1      (S1,S2) BIN FIXED (31) ,
9      1      (C1,S2) BIN FIXED (31) ,
10     1      (C1(4),C2(4)) DEC FIXED(31,5) ,
11     1      (L,LK,J,LG,N) BIN FIXED(31) ,
12     1      /PROCEDURE DE RECHERCHE DES VOISINS
13     1      D'UN CENTROIDE
14     1      VOISINS ; PROC(IND,CENT,YP1,YP2) ;
15     2      DCL IND BIN FIXED(31) ,
16     2      (YP1(4) ,YP2(4)) BIN FIXED(31) ,
17     2      CENT(4) DEC FIXED(31,5)
18     2      (D,IND1,IND2, DMIN1, DMIN2) BIN FIXED,
19     2      ;/FONCTION: DISTANCE ENTRE DEUX VECTEURS/
20     2      DIST : PROC(A,B) RETURN(BIN FIXED(31))
21     3      DCL(A(4),(4)) BIN FIXED(31);
22     3      DCL(I,SOMM) BIN FIXED (31);
23     3      SOMM=0 ;
24     3      DO I= 1 TO 4 ;
25     4      SOMM = SOMM+(5-I) * ((A(I)-B(I))**2
26     4      END ;
27     3      RETURN(SOMM) ;
28     3      END DIST ;
29     3      DMIN1=DIST(X(S(LG,IND) ),CENT) ;
30     2      YP1=X(S(LG,IND) , )

```

SUITE

```

31 2   DMIN2=DIST(X(S(LG,IND)+1,*)CENT) ;
32 2   YP2=X(S(LG,IND)+1,*) ;
33 2   IND1=S(LG,IND) ;
34 2   IND2=S(LG,IND+1)-1 ;
35 2   DO J=IND1 TO IND2 ;
36 3   D=DIST(X(J,*)CENT) ;
37 3   IF D < DMIN1 THEN DO ;
38 4   YP1=X(J,*) ;
39 4   DMIN1=D
40 4   END ;
41 3   IF (D < DMIN1 ET D < DMIN2) THEN DO ;
42 4   DMIN2=D ;
43 4   YP2=X(J,*) ;
44 4   END ;
45 3   END ;
46 2   END VOISINS ;

47 1   / PROCEDURE DE CLASSEMENT /
48 1   CLASSES : PROC(IND,YP1,YP2,BL1,BL2) ;
49 2   DCL (Q,IND1,IND2,J,IND,P,BL1,BL2,SIZE) BIN FIXED (31)
50 2   (YP1(4),YP2(4)) BIN FIXED (31)
51 2   INTERX(SIZE,4) BIN FIXED(31) CTL ;
52 2   /FONCTIN DISTANCE ENTRE DEUX VECTEURS/
53 2   DIST : PROC(A,B) RETURNS(BIN FIXED(31)) ;
54 3   DCL (A(4),B(4)) BIN FIXED(31) ;
55 3   DCL (I,SOMM) BIN FIXED(31) ;
56 3   SOMM=0 ;
57 3   DO I=1 TO 4 ;
58 4   SOMM=SOMM+(5-I)*((A(I)-B(I))*2) ;
59 4   END ;

```


SUITE

```

60 3 RETURN(SOMM) ;
61 3 END DIST ;
62 2 SIZE S(LG,IND+1)-S(LG,IND) ;
63 2 ALLOCATE INTERX ;
64 2 Q=1 ;
65 2 P=0 ;
66 2 IND1=S(LG,IND) ;
67 2 IND2=S(LG,IND+1)-1 ;
68 2 DO J=IND1 TO IND2 ;
69 3 IF DIST(X(J,*) ,YP1) =DIST(X(J,*) ,YP2) THEN DO ;
70 4 INTERX(Q,*)=X(J,*) ;
71 4 Q=Q+1 ;
72 4 END ;
73 3 ELSE DO ;
74 4 INTERX(S(LG,IND+1)-J(LG,IND)-P,*)=X(J,*) ;
75 4 P=P+1 ;
76 4 END ;
77 3 END ;
78 2
79 2 BL2=Q-1+S(LG,IND) ;
80 2 BL1=S(LG,IND) ;
81 2 Q=1 ;
82 2 DO J=IND1 TO IND2 ;
83 3 X(J,*)=INTERX(Q,*) ;
84 3 Q=Q+1 ;
85 3 END ;
86 2 FREE INTERX ;
87 2 END CLAES ;

```

88	1	CENTROIDE : PROC(LR, IND, CENT) ;
89	2	DCL (I, LR, IND1, IND2, J, IND, SOM) BIN FIXED(31) ;
90	2	CENT(4) DEC FIXED(31, 5) ;
91	3	DO J=1 TO 4 ;
92	3	SOM= 0 ;
93	3	IND1=S(LG, IND) ;
94	3	IND2=S(LG, IND+1)-1 ;
95	3	DO I=IND1 TO IND2 ;
96	4	SOM=SOM+X(I, J) ;
97	4	END (;
98	3	ENDY(J, SOM/(LG, IND+1)-S(LG, IND)) ;
99	3	END CENTROIDE ;
100	1	/ PROGRAMME PRINCIPALE /
101	1	
102	1	GET LIST(N) ;
103	1	ALLOCATE X
104	1	DO I=1 TO N ;
105	2	DO J=1 TO 4 ;
106	3	GET LIST(X(I, J)) ;
107	3	END ;
108	2	PUT EDIT X(I, %) (SKIP, 4(X(10), F(8))) ;
109	2	END ;
110	1	DO I=1 TO 6
111	2	DO J=1 TO 33 ;
112	3	S(I, J)=M+1
113	3	END ;
114	3	END ;
115	1	S(1, 1)=1
116	1	LG=1 ;
117	1	I=1 ;
118	1	CALL CENTROIDE (LG, I, C) ;
119	1	PUT EDIT (' ') (SKIP X(10) A(10))
120	1	Y(1, :) = 1 ;
121	1	DO LG=1 TO 5 ;

SUITE

```

122      2      DO I=1 TO 5 ;
123      3      C=Y((2 (LG-1))+I-1,*) ;
124      3      CALL VOISINS(I,C,Y1,Y2) ;
125      3      CALL CLASSES(I,Y1,Y2,S1,S2) ;
126      3      S(LG+1,2 I-1)=S1 ;
127      3      S(LG+1,2 I)=S2 ;
128      3      PUT EDIT(S1,S2)(SKIP,X(10),F(5),F(5)) ;
129      3      S1=2 I-1 ;
130      3      S2=2 I
131      3      PUT EDIT(' ')(SKIP X(10),A(11)) ;
132      3      CALL CENTROIDE(LK,S1,C1) ;
133      3      CALL CENTROIDE LK,S2,C2) ;
134      3      Y((2 LG)+(2 I)-2,*)=C1 ;
135      3      Y((2 LG)+(2 I)-1,*)=C2 ;
136      3      END ;
137      2      END ;
138      2      DO I=1 TO 33
139      1      DO J=S(6,I) TO S(6,I+1)-1) ;
140      2      PUT EDIT X(J,*) (SKIP,4(X(10),F(8))) ;
141      3      END ;
142      2      END ;
143      1      DO I=1 TO 63
144      2      PUT EDIT (Y(I,*))(SKIP,4(X(10),F(8,4))) ;
145      2      END ;
146      1      END NUUES

```

F I N D U P R O G R A M M E

SUITE

```

122      2      DO I=1 TO 5 ;
123      3      C=Y((2 (LG-1))+I-1,*) ;
124      3      CALL VOISINS(I,C,Y1,Y2) ;
125      3      CALL GLASSES(I,Y1,Y2,S1,S2) ;
126      3      S(LG+1,2 I-1)=S1 ;
127      3      S(LG+1,2 I)=S2 ;
128      3      PUT EDIT(S1,S2)(SKIP,X(10),F(5),F(5)) ;
129      3      S1=2 I-1 ;
130      3      S2=2 I
131      3      PUT EDIT(      )(SKIP X(10),A(11)) ;
132      3      CALL CENTROIDE(LK,S1,C1) ;
133      3      CALL CENTROIDE LK,S2,C2) ;
134      3      Y((2 LG)+(2 I)-2,*)=C1 ;
135      3      Y((2 LG)+(2 I)-1,*)=C2 ;
136      3      END ;
137      2      END ;
138      2      DO I=1 TO 33
139      1      DO J=S(6,I) TO S(6,I+1)-1) ;
140      2      PUT EDIT X(J,*) (SKIP,4(X(10),F(8))) ;
141      3      END ;
142      2      END ;
143      1      DO I=1 TO 63
144      2      PUT EDIT (Y(I,*)) (SKIP,4(X(10),F(8,4))) ;
145      2      END ;
146      1      END NUDES

```

F I N D U P R O G R A M M E

PROGRAMME DE RECONNAISSANCE (PHONE)

```

146      3      PHONE: PROC(V,ETIQ,D) ,
147      DCL V(4) BIN FIXED(31)
148      2      (D,D1,D2) BIN FIXED(31)
149      2      ETIQ CHAR(2) VARYING ,
150      2      TRANSL(32) CHAR(2) VARYING INIT(A1,A2,A5,A4,
151      2      A5,A7,A8,A9,B1,B2,B3,B4,B5,B6,B7,B8,B9,C1,
152      2      C2,C3,C4,C5,C6,C7,C8,C9,D1,D2,D3,D4,D5)
153      2      DCL(IND,IND1,IND2(K) BIN FIXED(31)
154      2      DIST: PROC(A,B) RETURNS(BIN FIXED(31))
155      3      DCL(A(4),(B(4) BIN FIXED(31)) ,
156      3      (I,SOM) BIN FIXED(31),H
157      3      SOM=0 ,
158      3      DO I=1 TO 4 ,
159      4      SOM SOM+(5-I) * (A(I) B I))**2
160      4      END ;
I61      3      RETURN(SOM) ;
I62      3      END DIST ;
I63      2      IND1=2 ,IND2=3 ;
I64      2      DO K=1 TO 4 ;
I65      3      D1=DIST(Y(IND1, ),V) ;
I66      3      D2 DIST(Y(IND2, ),V) ;
I67      3      IF D1 D2 THEN DO ;
I 68      4      IND=IND1 ;
I69      4      D=D1 ;
I70      4      IND1=2 IND1 ;
I7      4      IND2=2 IND1+1 ;
I72      4      END ;
I 73      4      END ;

```

```

I74      2      ETIQ=TRANSL(IND-3I) ;
I75      2      END PHONE ;
I76      I      / PROGRAMME PRINCIPAL /
I77      I      CALL NUÉE ;
I78      I      GET LIST(PRON) ;
I78      L      CALL PHONE(PRON,PHNNEM,DIST) ;
I79      I      PUT EDIT (PRON)(SKIP,4(X(IO),F(5))) ;
I80      I      PUT EDIT('PHONÈME RECONNU')(SKIP X(IO),A(35)) ;
I81      I      PUT EDIT (PHONÈME)(SKIP,X(IO),A(2)) ;
I82      I      PUT EDIT ('XXXXXXXXXXXXXXXXXXXX')(SKIP,X(IO),A(36)) ;
I83      I      END TEST ;

```


Resultats obtenus

Nous avons élaboré un programme constitué essentiellement de deux parties:

- un programme d'apprentissage (NUEE)
- " " " " de reconnaissance (PHONE)

Nous avons simulé 64 données en entrée, le programme élabore un dictionnaire de 32 classes. Parmi ces 32 classes obtenues, 16 sont vides. Cela s'explique par le fait que la taille de l'échantillon simulé à l'entrée est réduite.

Ensuite nous avons ^{simulé} une donnée (un vecteur) en vue de sa reconnaissance. Le programme de reconnaissance ne fournit pas de résultat car il dépend du programme d'apprentissage.

CONCLUSION:

La reconnaissance automatique de la parole est un sujet difficile et complexe, qui est toujours au stade de la recherche fondamentale.

La simple analyse acoustique du message parlé ne peut suffire à le résoudre. La parole n'est pas une simple juxtaposition de phonèmes pour sa reconnaissance il faut faire intervenir, en plus des notions de critère et de phonème, la notion de syllabe et de la redondance liée au vocabulaire et aux idées.

On a vu qu'à travers ce travail les différentes étapes qui interviennent dans un problème de reconnaissance de la parole peuvent être abordés par diverses techniques mathématiques, leur mise en oeuvre demande beaucoup de simplification et d'approximation mais elles permettent de fonder la reconnaissance sur des bases solides. Le choix de la méthode d'analyse dépend du type de méthode de reconnaissance utilisée et du résultat qu'on veut obtenir.

On distingue principalement deux méthodes de reconnaissance l'une dite globale basée sur une comparaison directe du signal à des descriptions de mots ou de séquence de mots en termes de modèles acoustiques. Cette méthode semble devoir aboutir dans un avenir proche à des réalisations opérationnelles. Ses limites tiennent surtout au volume de calcul et de mémoire, croissant directement avec la taille du vocabulaire utilisable.

La deuxième méthode est basée sur la reconnaissance des phonèmes, elle comporte des étapes: le prétraitement, la reconnaissance des phonèmes et le décodage lexical. Les performances des procédures de reconnaissance des phonèmes sont médiocres et dépassent rarement 60%; ce qui complique davantage la tâche du décodage lexical. L'évaluation des résultats de la segmentation et d'identification n'est pas aisée car plusieurs facteurs entrent en jeu. En ce qui concerne la segmentation en phonèmes, on peut dire que les résultats sont supérieurs à 90% avec en général peu de phonèmes omis (2 à 10 %) et parfois plus de phonèmes insérés (5 à 15 %). Les résultats de l'identification des phonèmes est difficile à évaluer, on peut seulement dire pour une trentaine de phonèmes prononcés par un ou quelques locuteurs, les résultats varient de 50 à 70% de reconnaissance.

B I B L I O G R A P H I E

- GUIBERT, J. " La parole: compréhension et synthèse par les ordinateurs" P.U.F(1979)
- LEINARD, J.S (1977) "Les processus de la communication parlée. Introduction à l'analyse et la synthèse "
- MAX, J. (1981) "Méthodes et techniques de traitement du signal et application aux mesures phonétiques" Tome I , 3^{ième} édition Masson .
- BENZEKRI, A "Analyse de données" Tome I taxonomie
- MERCIER, G. (1977) "Analyse acoustique et transcription phonétique du signal de la parole" Information générale IG/DAS/SST/I
- BOUSSEKSOU, B. (1983) " Reconnaissance automatique de la parole par les méthodes globales. Application aux particularités linguistiques de l'arabe standard" Thèse de magister
- GUERTI, . (1984) "Contribution de la synthèse de la parole en arabe standard" Thèse de magister
- GRENIER, A. (1977) "Identification du locuteur et adaptation au locuteur d'un système de RECONNAISSANCE phonémique" Thèse de docteur-ingénieur.