

وزارة الجامعات  
Ministère aux Universités

ECOLE NATIONALE POLYTECHNIQUE

DEPARTEMENT électronique

المدرسة الوطنية المتعددة التقنيات  
BIBLIOTHEQUE — المكتبة  
Ecole Nationale Polytechnique

PROJET DE FIN D'ETUDES

SUJET

RECONNAISSANCE  
AUTOMATIQUE DE LA  
PAROLE PAR LA  
METHODE STOCHASTIQUE

Proposé par : B. Bousseksou

Etudié par : Z. Ait-kheddache

Dirigé par : B. Bousseksou

PROMOTION 91 = 92

<< AU NOM D'ALLAH >>

سجلت في المكتبة  
في تاريخ 1/10/2005  
على رقم 1000  
Les documents et les  
références

### REMERCEMENTS

Je tiens <sup>S</sup> à remercier tous ceux qui m'ont aidé et encouragé durant mes dernières années d'étude.

Je remercie vivement Mr Bousseksou, mon promoteur, qui a su me guider et me conseiller pour mener ce projet à son terme.

Je remercie, de même, le personnel de la Bibliothèque et du centre de calcul, pour leur amabilité et toute l'aide précieuse et irremplaçable qu'ils m'ont fournie.

**TABLE DES MATIERES**

<b>TABLE DES MATIERES</b> .....	1
<b>BIBLIOGRAPHIE</b> .....	3
<b>INTRODUCTION GENERALE</b> .....	4
<b>CHAPITRE I LA PAROLE NATURELLE.</b> .....	5
I.1 - INTRODUCTION .....	5
I.2 - PRODUCTION DE LA PAROLE NATURELLE.....	5
I.3 - MODELISATION DE LA PRODUCTION DE LA PAROLE .....	7
I.4 - LES PHONEMES .....	10
I.5 - CONCLUSION .....	12
<b>CHAPITRE II ANALYSE DE LA PAROLE</b> .....	13
II.1 - INTRODUCTION .....	13
II.2 - PRETRAITEMENT .....	13
II.2.1 - ECHANTILLONAGE .....	13
II.2.2 - PREACCENTUATION .....	14
II.2.3 - FENETRAGE.....	14
II.3 - ANALYSE CEPSTRALE .....	16
II.4 -ANALYSE PAR PREDICTION LINEAIRE .....	20
II.5 -CONCLUSION .....	26
<b>CHAPITRE III RECONNAISSANCE DE LA PAROLE</b> .....	27
III.1 - METHODES STOCHASTIQUES .....	28
III.2 - SEGMENTATION DE LA PAROLE .....	29
III.3 - MODELISATION DES PHONEMES EN SEGMENTS STOCHASTIQUES. 31	
III.3.1 - LA NORMALISATION TEMPORELLE.....	31
III.3.2 - LE MODELE PROBABILISTE.....	33
III.4 - ALGORITHMES DE RECONNAISSANCE.....	34
III.5 - LA PHASE D'APRENTISSAGE.....	42
III.5.1 - ESTIMATION DES PARAMETRES.....	43

III.5.2 - SEGMENTATION AUTOMATIQUE.....	44
CONCLUSION GENERALE.....	47
ANNEXES :	
A - NOTIONS DE PROBABILITES.....	48
B - NOTIONS DE STATISTIQUES.....	51
C - RESULTAS DE LA SIMULATION.....	52

## BIBLIOGRAPHIE

- [1] R. BOITE : Traitement de la Parole
- [2] F. CINARE : Reconnaissance et Synthèse de Parole
- [3] IEEE : A Stochastic Segment Model for Phoneme - Based  
Continuous Speech Recognition (Décembre 1989)
- [4] IEEE : LPC Speech Coding Based On Variable-Length Segment  
Quantization ( Septembre 1988 )
- [5] E. EMERIT: Cours de Phonétique Acoustique
- [6] J. LIFITMAN: Les Méthodes Rapides de Transformation du  
Signal
- [7] A. SPATARU: Fondements de La Théorie de Transmission de  
l'Information

## INTRODUCTION GENERALE



Les machines qui parlent et qui reconnaissent la parole vont avoir beaucoup d'applications utiles dans le futur. Beaucoup de temps pourrait-être économisé, en utilisant un ensemble d'instructions parlées à la place d'un ensemble d'instructions imprimées, quand il faudra mettre en marche un équipement complexe, dans les situations où les yeux et les mains sont déjà occupés, où des actions peuvent-être représentées sous formes d'informations nouvelles qui peuvent être mieux représentées sous forme de messages parlés.

Déjà, on peut trouver sur le marché des produits aux capacités limitées certes, mais pouvant rendre de grands services aux utilisateurs potentiels de cette technologie.

En reconnaissance de la parole continue, dans les vocabulaires larges, les mots sont fréquemment modélisés comme étant des suites d'unités linguistiques telles que les phonèmes, autrement dit, un mot est modélisé acoustiquement en concaténant les modèles acoustiques des phonèmes suivant sa chaîne structurée de prononciation entraînée dans le dictionnaire contenant les "transcriptions" des différents phonèmes. L'Avantage de cette approche est qu'il n'est pas nécessaire d'entraîner tous les mots dans le vocabulaire mais uniquement les modèles des différents phonèmes.

La méthode des chaînes de Markov est une modélisation probabiliste de la réalisation acoustique d'un phonème, mais bien que cette dernière a donné de bons résultats en ce qui concerne la modélisation des phonèmes aux longueurs variables, la reconnaissance de la parole est loin d'atteindre la performance, chose pour laquelle nous proposons à travers ce modeste travail, une nouvelle approche appelée MODELISATION DES PHONEMES EN SEGMENTS STOCHASTIQUES dans laquelle chacun des segments phonétiques, aux longueurs différentes est transformé en un segment de longueur fixée commune.

## CHAPITRE I

### LA PAROLE NATURELLE

#### I - 1/- INTRODUCTION

Le phénomène de la production de la parole suscite aujourd'hui l'intérêt des chercheurs. En effet, une meilleure connaissance de ce phénomène est devenue nécessaire pour pouvoir avancer dans les secteurs de la reconnaissance de la parole. Une telle connaissance exige l'étude de plus en plus approfondie des caractéristiques acoustiques de la parole; ainsi que celle des différentes divisions possibles de la parole.

#### I-2/ - PRODUCTION DE LA PAROLE NATURELLE

Le discours de la parole est une succession de différences de pression de l'air qui engendrées par le système de la phonation (fig I-1)

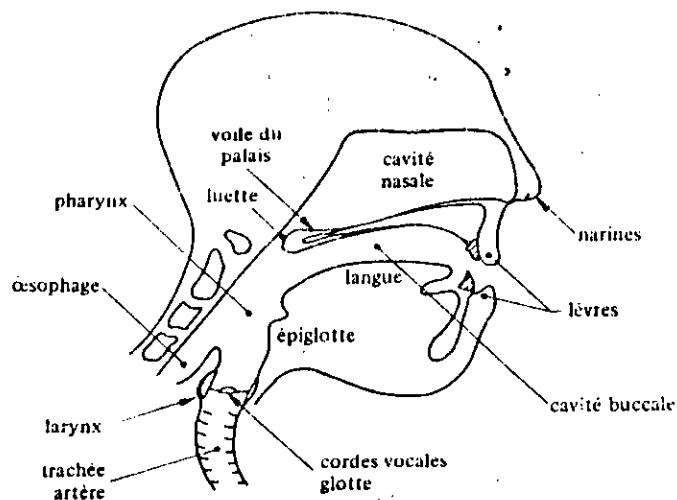


Fig I.1 Système phonatoire humain

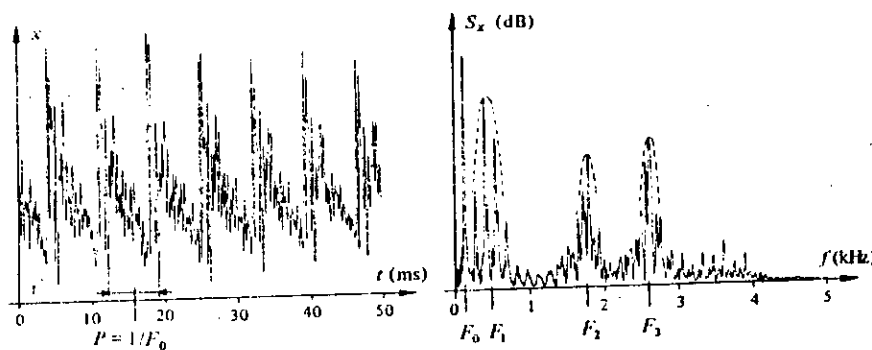


Le conduit vocal humain peut être considéré comme une succession de tubes ou cavités acoustiques de réactions diverses.

Les sons voisés résultent donc de l'excitation du conduit vocal par des impulsions périodiques de pression liées aux oscillations des cordes vocales: l'ouverture brusque de la glotte libère la pression accumulée en amont; elle se referme ensuite plus graduellement.

Un son voisé est un signal quasi-périodique (fig I-2 et I-3).

Sur la figure I-3, on observe les raies qui correspondent aux harmoniques du fondamental  $F_0$ . L'enveloppe de ces raies présente des maximums appelés formants et qui correspondent aux fréquences propres  $F_i$  ( $i=1, 2, 3, \dots$ ) du conduit vocal.



I-2 et I-3: Un signal vocal et son spectre

Les trois premiers formants sont essentiels pour caractériser le spectre vocal, les formants d'ordre supérieur ont une influence limitée. Un son non voisé ne présente pas de structure périodique, il peut être considéré comme un bruit blanc filtré par la transmittance du conduit vocal. Remarquons toutefois qu'un son non voisé présente la même structure formantique que le même son voisé.

### I-3/- Modélisation de la production de la parole

La modélisation proposée ici utilise le formalisme des systèmes échantillonnés (Transformée en Z).

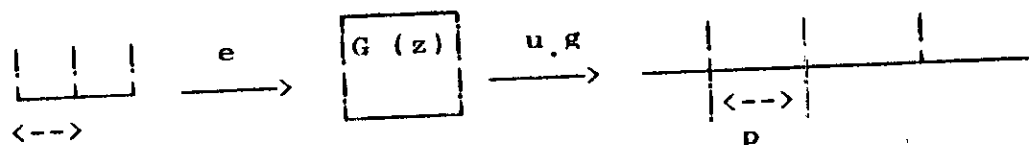


fig I-4 Modélisation de la source pour les sons voisés

Pour les sons voisés, la source est un train périodique d'ondes (fig 4). Ce train d'ondes est modélisé par un filtre passe-bas d'ordre 2 à pôles réels dont la fréquence de coupure est de l'ordre de 100 HZ:

$$G(z) = \frac{A}{(1+\alpha z^{-1})(1+\beta z^{-1})}$$

Pour les sons non voisés, la source est un bruit blanc.

On peut assimiler le conduit vocal à une cascade de résonateurs dont la transmittance est de la forme:

$$V(z) = \frac{B}{\prod_{k=1}^M (1 + b_{1k} z^{-1} + b_{2k} z^{-2})}$$

Chaque résonateur correspond à un formant dont la fréquence centrale est donnée par :

$$F_k = (1/2\pi) \cdot f_0 \cdot \cos^{-1} \left( \frac{-b_{1k}}{2 \cdot \sqrt{b_{2k}}} \right)$$

$f_0$  : fréquence d'échantillonnage.

Le son est finalement émis à travers l'ouverture des lèvres, où :

$$R(z) = C (1 - z^{-1})$$

En résumé, la transmittance globale entre le train d'impulsions et le son émis :

$$T(z) = G(z) \cdot V(z) \cdot R(z)$$

$$= \frac{\sigma \cdot (1 - z^{-1})}{(1 + \alpha \cdot z^{-1}) \cdot (1 + \beta \cdot z^{-1}) \cdot \prod_{k=1}^M (1 + b_{1k} \cdot z^{-1} + b_{2k} \cdot z^{-2})}$$

On suppose que l'un des pôles de  $G(z)$  est proche de l'unité donc:

$$T(z) = \frac{\sigma}{(1+\alpha z^{-1}) \prod_{k=1}^M (1 + b_{1k} z^{-1} + b_{2k} z^{-2})} = \frac{\sigma}{A(z)}$$

On a posé: 
$$A(z) = (1 + \alpha z^{-1}) \prod_{k=1}^M (1 + b_{1k} z^{-1} + b_{2k} z^{-2})$$

$$= 1 + \sum_{i=1}^{2M+1} (a_i z^{-i})$$

La transmittance de ce modèle est dite tout-pôles; son inverse le polynôme  $A(z)$  est la transmittance du filtre inverse. Ses limitations sont cependant évidentes. En premier lieu, la source est soit un train périodique d'impulsions, soit un bruit blanc; les sons fricatifs voisés ( $v, z, \dots$ ) ne peuvent pas être produits par ce modèle.

En second lieu, la production des sons nasalisés fait intervenir deux cavités associées en parallèle; la transmittance correspondante est de la forme:

$$\frac{\sigma_1}{A_1(z)} + \frac{\sigma_2}{A_2(z)} = \frac{\sigma_1 A_2(z) + \sigma_2 A_1(z)}{A_1(z) A_2(z)}$$

Toutefois, malgré ses limitations, la transmittance tous-pôles est la base de la modélisation par prediction linéaire. Pour pouvoir assimiler le numérateur à une constante, on doit surestimer le degré du dénominateur.

On remarque aussi que le signal vocal est stationnaire pendant des intervalles de temps de l'ordre de 20 ms. Durant cet intervalle, les coefficients de  $T(z)$  sont constants.

#### I-4/ - LES PHONEMES

Les linguistes ont défini le phonème comme étant la plus petite unité phonétique. La langue française contient 36 phonèmes comprenant 16 voyelles, 17 consonnes, et 3 semi-consonnes. Sur la figure I.5, on peut voir une classification des phonèmes.

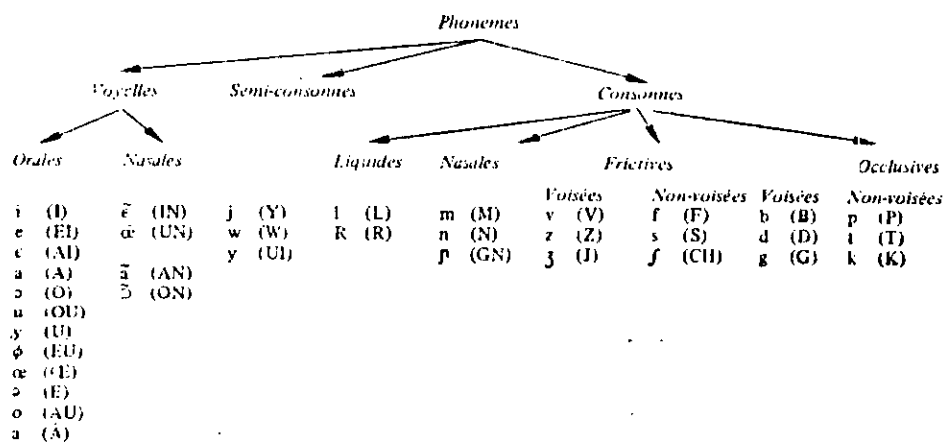


fig I-5 phonèmes de la langue française, d'après Kunt [9]

La production d'un phonème donné laisse toutefois place à une certaine variabilité sur le plan acoustique.

Les voyelles orales (i, o, u, ...) sont émises sans intervention de la cavité nasale. Pour les voyelles nasales, le conduit nasal est couplé à la cavité buccale et l'émission se produit à la fois par les narines et par la bouche. Les nasales font aussi intervenir la cavité nasale.

Les consonnes fricatives résultent de l'écoulement de l'air dans une constriction étroite en un point du conduit vocal, en particulier au niveau des lèvres et des dents. Les sons fricatifs sont non voisés (f, s, ...) ou voisés (v, z, ...).

Les consonnes occlusives correspondent à des sons essentiellement dynamiques. Une forte pression est créée en amont d'une occlusion maintenue en un certain point du conduit vocal puis relâchée brusquement. La période d'occlusion est appelée la phase de tenue.

Pour les occlusives voisées un son basse fréquence est émis par vibration des cordes vocales pendant la phase de tenue; pour les occlusives non voisées, la tenue est un silence.

	1 <sup>er</sup> for. (Hz)	2 <sup>ème</sup> for. (Hz)
a	750	1350
o	375	750
i	250	2500
e	400	2200
u	250	600

Fig I-6 les deux premiers formants des voyelles

La table I-6 représente les voyelles pour les deux premiers formants. Les deux premiers formants caractérisent déjà beaucoup mieux toutes les voyelles.

### I.5/ - CONCLUSION

Le mécanisme de la production de la parole qui a été décrit dans ce chapitre conduit très naturellement à une modélisation particulière appelée MODELISATION AUTOREGRESSIVE. Nous disposons là d'un outil très commode pour procéder à l'analyse du signal vocal. On prendra quand-même quelques précautions car l'hypothèse de stationnarité du signal n'est valable que durant de courts intervalles de temps. Les procédures de modélisation doivent être adaptées pour en tenir compte.

CHAPITRE II



## CHAPITRE II

### ANALYSE DE LA PAROLE

#### II-1 INTRODUCTION :

Pour pouvoir réaliser un système de reconnaissance de la parole, il est très important de faire passer ce signal par une analyse qui consiste à tirer du signal les paramètres pertinents capable de représenter correctement toutes ses caractéristiques.

La première opération d'analyse consiste à tronçonner le temps en intervalles d'environ 10 à 20 ms. Ces intervalles sont appelés "fenêtres temporelles". L'analyse du signal parole repose sur l'hypothèse fondamentale suivante :

" Le signal à analyser est stationnaire pendant toute la durée de la fenêtre d'analyse, c'est à dire que ses propriétés en particulier spectrales ne varient pas sur la longueur de la fenêtre ".

L'analyse sera faite successivement sur chacune des fenêtres. Parmi les méthodes d'analyse citons :

L'analyse cepstrale; l'analyse par prédiction linéaire et l'analyse par banc de filtre; peut être numérique (logiciel) ou matériel (vocodeur à canaux par exemple).

#### II-2 PRETRAITEMENT

Avant d'analyser le signal parole, on doit passer par un prétraitement qui consiste en un échantillonnage, une préaccentuation et un fenêtrage.

##### II-2-1 ECHANTILLONNAGE

Dans le cas des méthodes numériques, ce traitement ne s'effectue pas directement sur les signaux analogiques à temps continu fournis par un microphone.

Les signaux analogiques sont échantillonnés, codés, puis rangés sous forme numérique dans une mémoire pour être traité par ordinateur.

Le spectre du signal parole est limité à 6 KHZ, il conserve donc ses caractéristiques.

Selon le théorème de SHANNON :  $F_e \geq 2 \cdot F_m$

Avec  $F_e$  : Fréquence d'échantillonnage

$F_m$  : Fréquence maximale du signal parole

donc en théorie 12 KHZ vont suffire comme fréquence d'échantillonnage, et si on prend 256 échantillons par fenêtre de 20 ms on aura pour fréquence d'échantillonnage :

$$F_e = \frac{256}{2 \times 10^{-2}} = 12.8 \text{ KHZ}$$

### II-2-2 PREACCENTUATION :

Du fait de l'évolution du spectre du signal parole en haute fréquence dans les deux milieux : Le conduit vocal et l'air extérieur, il subit une baisse d'énergie de 6 dB par octave.

La préaccentuation est nécessaire pour rétablir le niveau énergétique . C'est le cas des sons voisés car le spectre est surtout localisé dans la partie haute fréquence.

### II-2-3 FENETRAGE :

Le fenêtrage limite l'amplitude des rebonds fréquentiels dus à la limitation temporelle du signal.

Il existe plusieurs types de fenêtres :

- Fenêtre rectangulaire :

$W_r(k) = A$  pour  $k = 0, N$  Avec  $A$  : constante

pour un nombre d'échantillons  $N = 9$ , L'atténuation = -13 dB .

- Fenêtre triangulaire :

$W(k) = 1 - 2|k| / N$  pour  $|k| \leq N / 2$

L'amplitude des lobes secondaires est inférieure à celle de la fenêtre rectangulaire, mais cette fenêtre ne nous paraît pas très performante du fait de l'élargissement du pic central.

- fenêtre de Hanning :

$$W_h(k) = 1 / 2 ( 1 + \cos ( 2 \pi \cdot k / N ) ) \quad \text{pour } |k| \leq N / 2$$

Son premier lobe secondaire n'est pas vraiment atténué, par contre les autres lobes secondaires sont atténués à environ -43 dB du pic central.

- Fenêtre de Hamming :

$$W_h(k) = 0,54 - 0,46 \cos ( 2 \cdot \pi \cdot (k - 1) / N )$$

pour  $k = 1, N$ . Elle a 99,96 % de son énergie dans le lobe principal et le lobe secondaire .

Le plus important est de 40 dB en dessous du lobe principal .

- Critère de choix de la fenêtre :

Le choix de la forme particulière d'une fonction fenêtre dépend principalement de la largeur du pic central et de l'amplitude des lobes secondaires relative à celle du pic central . Si le pic central est large, les transitions rapides de la transformée de Fourier du signal original sont très mal approximées.

D'après le résultat selon lequel 99,96 % d'énergie est concentrée dans le lobe principal, le niveau du premier lobe secondaire est 43,9 au dessous du lobe principal; toutes ces qualités nous conduit à opter pour la fenêtre de Hamming.

### II-3 ANALYSE CEPSTRALE :

Le signal de parole contient des informations phonétiques essentiellement contenues dans le spectre du signal et des informations prosodiques essentiellement contenues dans le pitch ( la fréquence du fondamental ).

Le conduit vocal module le signal pulsé qui est fourni par la vibration des cordes vocales.

Il y a combinaison par convolution des deux filtres qui sont les cordes vocales et le conduit vocal.

soit  $h(n)$  le signal issu de la source d'excitation et  $x(n)$  la fonction de transfert du conduit vocal indépendante de  $h(n)$

$$Y(n) = \sum_{k=-\infty}^{k=+\infty} [h(n-k) \cdot x(k)] = h(n) * x(n) \quad (\text{II-1})$$

si  $D$  est l'opérateur de déconvolution

$$D[Y(n)] = D[h(n) * x(n)] = D[h(n)] + D[x(n)] \quad (\text{II-2})$$

$$Y'(n) = h'(n) + x'(n)$$

$h(n)$  et  $x(n)$  sont régulièrement échantillonnés,  $*$  est l'opération de convolution .

Par la transformée en  $Z$  on obtient :

$$Y(Z) = H(Z) \cdot X(Z)$$

$$\ln|Y(Z)| = \ln|H(Z)| + \ln|X(Z)| \quad (\text{II-3})$$

Le cepstre " $C$ " sera alors la transformée inverse de  $Z$  de

$\ln|Y(n)|$  " telque :

$$C(n) = (1/N) \cdot \sum_{k=1}^N (\text{Log} |Y(f)| e^{i2\pi n f / N}) \quad (\text{II-4})$$

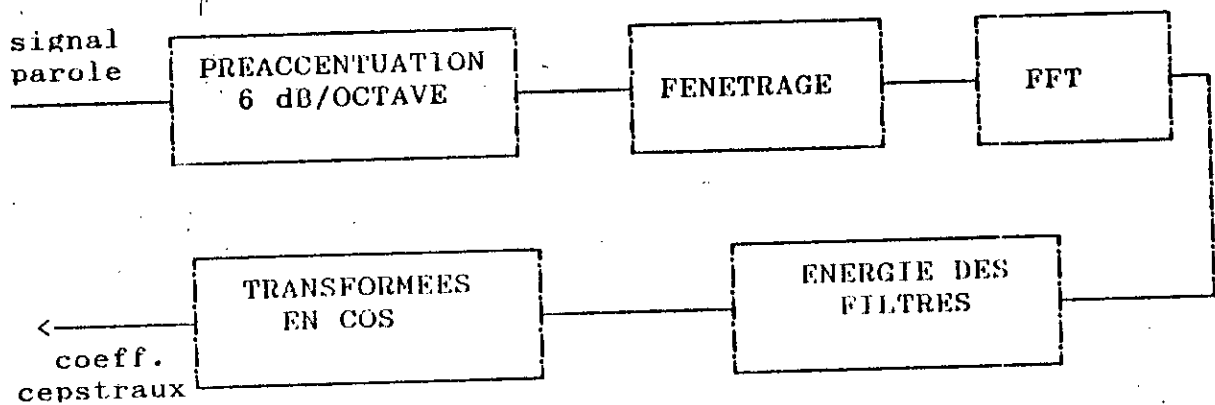
avec N: nombre de points dans une fenêtre  
 n: rang du coefficient cepstral

Les coefficients cepstraux se calculent par l'équation :

$$C(n) = (1 / N_f) \cdot \sum_{k=1}^{N_f} \text{Log} [|E(k)| \cos[n (k-1) / N_f]] \quad (\text{II-5})$$

avec: n: le rang du coefficient cepstral  
 k: le numéro du filtre d'énergie E(k)  
 N<sub>f</sub>: le nombre total des (16 à 25) filtres triangulaires

ETAPES D'ANALYSE :



ETAPES DE L'ANALYSE

La préaccentuation est une dérivation numérique :

$$Y(k+1) = Y(k+1) - Y(k)$$

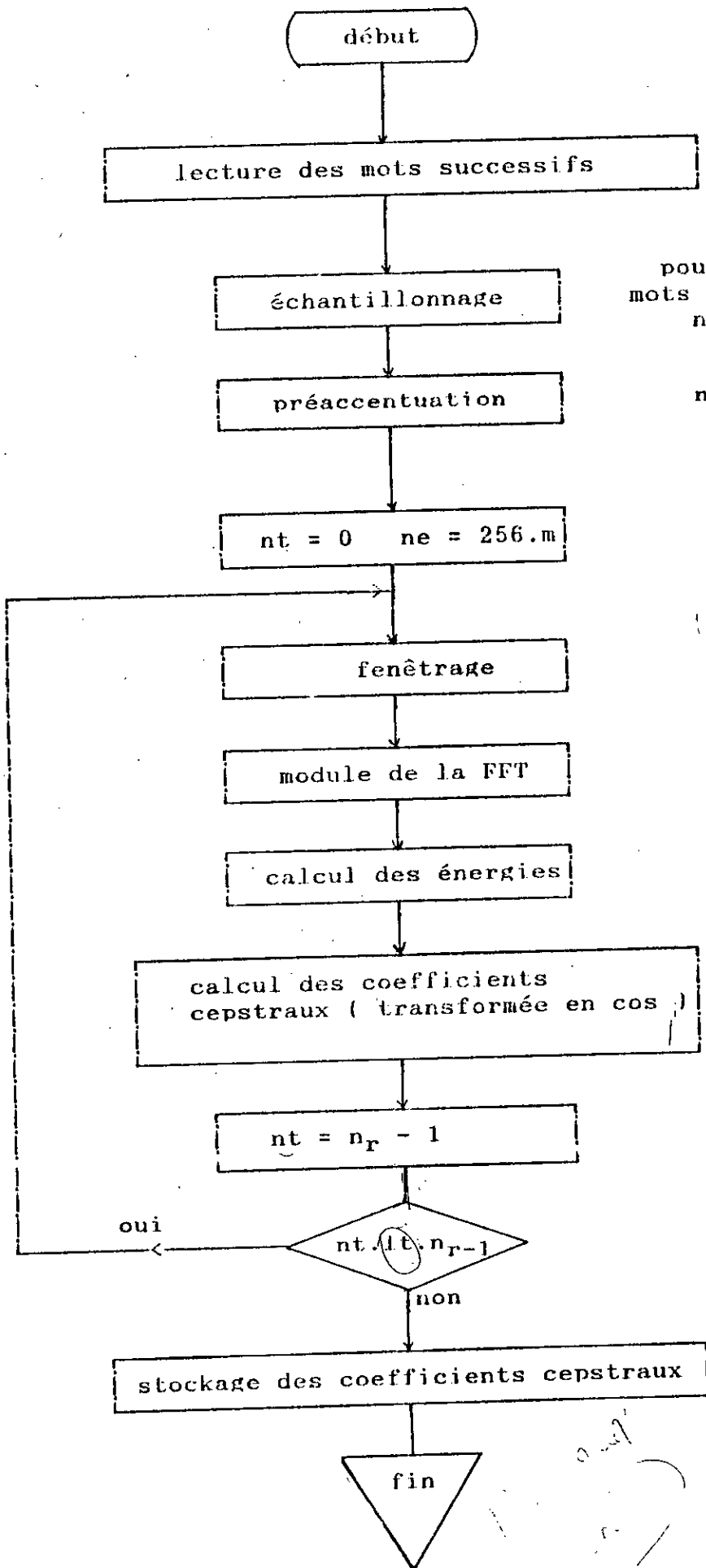
Le calcul du spectre se fait grâce à la FFT.

- Echelle Mel : l'oreille est sensible à une échelle quasi-logarithmique dite échelle Mel qui est linéaire sur le premier KHZ et logarithmique au-delà. L'échelle logarithmique tente d'éviter de donner un même poids à des zones de fréquence qui n'ont pas la même densité d'information.

### ORGANIGRAMME

Nous présentons ici l'organigramme de la méthode d'analyse cepstrale.

si on prend par exemple un phonème de 0,8 seconde, celui-ci sera tronçonné en 40 trames de 20 ms chacune et dans chaque trame nous avons 256 échantillons.



pour chacun des mots successifs on a  
 nt: nombre de trames  
 ne: nombre d'échantillons  
 nr: nombre total d'échantillons du mot

## II-4/ ANALYSE PAR PREDICTION LINEAIRE :

Cette technique permet d'estimer des paramètres comme le pitch, les formants, le spectre instantané.

Le signal est approximé par un polynôme d'ordre p.

Un échantillon de signal parole est représenté par les p échantillons qui le précèdent par une combinaison linéaire.

$$S_p(n) = \sum_{k=1}^p [a(k) \cdot S(n-k)] \quad (\text{II-6})$$

n = 1, N

p: ordre de prédiction

Le signal parole n'étant pas polynômial, les coefficients LPC a(k) sont calculés en minimisant la somme des carrés des différences entre les échantillons réels de la parole s(n) et la valeur estimée par combinaison linéaire S<sub>p</sub>.

$$E = \sum_{n=1}^{N+p} [ [ S(n) - \sum_{k=1}^p a(k) \cdot S(n-k) ]^2 ] \quad (\text{II-7})$$

Il existe plusieurs méthodes de minimisation de l'erreur quadratique citons entre autres :



- La Méthode de covariance
- La Méthode d'autocorrélation
- La Méthode en treillis

La prédiction linéaire est une méthode précise d'estimation des paramètres de la parole de plus, les calculs de certaines méthodes de résolution sont relativement rapides.

Dans notre cas, la sommation se fait sur un nombre N fini d'échantillons pendant lesquels les caractéristiques du conduit vocal et de la source sont constantes, ce qui conduit à la stationnarité de l'échantillon considéré, pour cela on utilise la méthode d'autocorrélation.

- Méthode d'autocorrélation :

Détermine les coefficients  $a(k)$  pour lesquels E est minimale.

Cela se fait en annulant la dérivée de E :

$$\frac{\partial E}{\partial R(i)} = \sum_{n=1}^{N+p} 2[s(n) - \sum_{k=1}^p a(k) \cdot s(n-k)] \cdot S(n-i) = 0 \quad (\text{II-8})$$

$$\sum_{n=1}^{N+p} S(n) \cdot S(n-1) = \sum_{k=1}^p a(k) \cdot \sum_{n=1}^{N+p} S(n-k) \cdot S(n-i) \quad (\text{II-9})$$

on pose :

$$c(i,k) = \sum_{n=1}^{N+p} S(n-k) \cdot S(n-i) \quad (\text{II-10})$$

$$\sum_{n=1}^{N+p} S(n).S(n-i) = \sum_{k=1}^p a(k).c(i,k) \quad (\text{II-11})$$

Avec  $i = 1, p$  et  $k = 1, p$

d'où :

$$c(i,k) = \sum_{n=1}^{N+i-k} S(n).S(n+i-k) \quad (\text{II-12})$$

$c(i,k)$  est la matrice d'autocorrelation. c'est une matrice carrée d'ordre "p" dite aussi "matrice de Toeplitz" car les éléments situés symétriquement de part et d'autre de la diagonale sont égaux.

on note  $R(k)$  la fonction d'autocorrelation définie par :

$$R(k) = \sum_{n=1}^{N+k} S(n).S(n+k) \quad (\text{II-13})$$

avec  $k=1, p$

$$R(k) = R(-k)$$

$$\text{d'où: } c(i,k) = R(|i-k|)$$

(II-14)

L'équation II-8 peut s'écrire sous forme :

$$\sum_{K=1}^p a(k).R(|i-k|) = R(i) \quad (\text{II-15})$$

avec  $i=1, p$

L'équation II-14 peut s'écrire sous forme matricielle

$$\begin{bmatrix} R(0) & R(1) & \dots & R(p-1) \\ R(1) & R(0) & \dots & R(p-2) \\ \vdots & \vdots & \ddots & \vdots \\ R(p-1) & R(p-2) & \dots & R(0) \end{bmatrix} \begin{bmatrix} a(1) \\ \vdots \\ a(p) \end{bmatrix} = \begin{bmatrix} R(1) \\ \vdots \\ R(p) \end{bmatrix} \quad (\text{II-16})$$

De nombreuses méthodes permettent de résoudre le système d'équation linéaire telles que :

- méthode de Gauss-Seidel
- Méthode de Jacobi
- Méthode de Gauss-Jordan
- Méthode de Durbin

Mais en tenant compte de la rapidité d'exécution et de l'encombrement mémoire réduit, nous avons choisi pour résoudre le système d'équation d'auto-correlation la méthode de Durbin.

Méthode de Durbin :

$$D(i) = R(i) - \sum_{j=1}^{i-1} [A(j, i-1) \cdot R(i-j) / E(j-1)] \quad (\text{II-17})$$

$$E(i) = [1 - D^2(i)] \cdot E(i-1) \quad (\text{II-18})$$

$$A(j, i) = A(j, i-1) - D(i) \cdot A(i-j, i-1)$$

$$A_1(j) = a(j, 12)$$

où

$D(i)$  : coefficients de réflexion

$E(i)$  : Erreur quadratique

$A(j, i)$  : coefficients de prédiction

$A_1(j)$  : coefficients d'autocorrelation

avec  $i=1, 12$  et  $j=1, i-1$

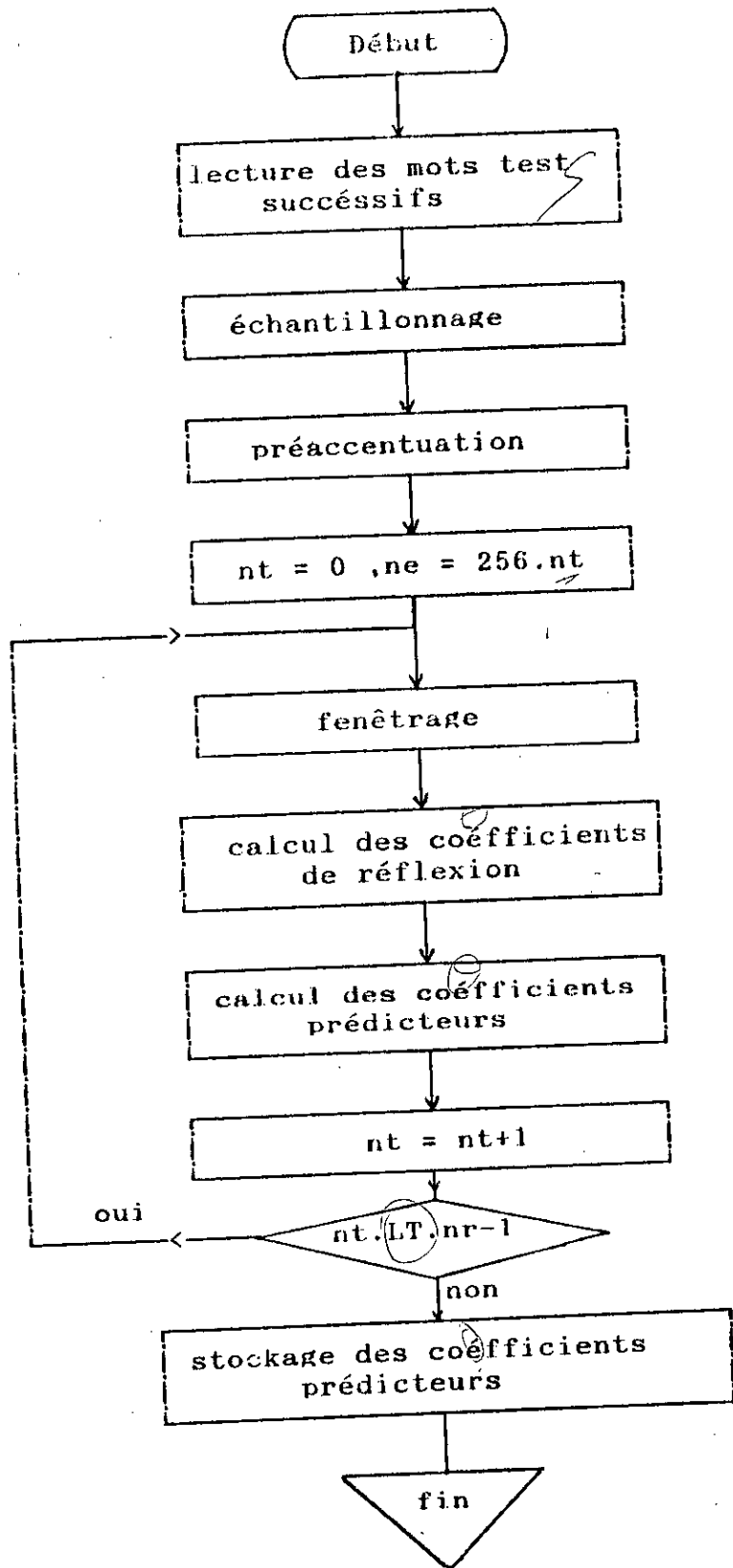
avec les conditions initiales suivantes :

$$E(1) = R(1)$$

$$D(2) = [R(2)/R(1)]$$

$$A(2, 2) = D(2)$$

$$E(2) = [1 - D^2(2)] \cdot R(1)$$



## II-5 CONCLUSION :

Dans ce chapitre nous avons étudié deux méthodes d'analyse " la méthode cepstrale " et " la méthode de prédiction linéaire ". La méthode cepstrale inverse (le cepstre) du signal et favorise les hautes fréquences par rapport aux basses fréquences qui ne présentent pas le même intérêt puisqu'elles sont spécifiques au locuteur (pitch).

De plus la méthode cepstrale est économique du point de vue représentation étant donné que 8 coefficients au lieu de 12 par trame (comme c'est le cas de la LPC) suffisent à bien représenter le signal.

Il est également possible grâce à la méthode cepstrale de réaliser une économie dans le temps de calcul.

Faire un traitement pour 8 coefficients est évidemment plus rapide que de le faire pour 12, ce qui nous rapproche de l'aspect temps réel qui est un des buts les plus recherchés.

**CHAPITRE III**

## CHAPITRE III

### RECONNAISSANCE DE LA PAROLE

Il a toujours été considéré que l'oreille procède à une certaine analyse spectrale continue dont les résultats sont transmis au cerveau. L'identification d'un mot consiste donc en une succession de comparaisons de spectres, si l'on se limite au niveau acoustique.

Un spectre peut être considéré d'une façon objective par différents ensembles de paramètres: en particulier le spectre du modèle AR, qui correspond très bien à l'enveloppe du spectre vocal, peut être caractérisé par les coefficients  $a_p(i)$ .

Un ensemble ordonné de paramètres caractérisent un spectre vocal est appelé vecteur acoustique ou vecteur spectral.

Un mot est bien sûr, constitué par une suite de vecteurs acoustiques.

La principale difficulté consiste à reconnaître un mot. Pour cela on dispose de trois approches.

> La première exploite la notion de distance entre deux mots.

L'algorithme DTW permet d'optimiser cette distance. Sur cette méthode, nous ne donnerons aucun détail. La deuxième approche apparemment très différente consiste à spéculer sur un modèle statistique de production pour chaque mot du vocabulaire. Et la troisième, qui ne diffère pas trop de la deuxième, consiste à modéliser les phonèmes en segments stochastiques.

Nous allons approfondir et détailler cette méthode.



### III.1 METHODES STOCHASTIQUES

On distingue deux méthodes:

La première est celle qui a déjà été citée et qui s'appuie sur l'hypothèse qu'un automate peut émettre un mot ou un phonème. A chaque état interne de l'automate correspond l'émission d'un vecteur acoustique.

C'est donc la succession des états internes qui provoque la production d'un mot. On dit de ce modèle qu'il est doublement stochastique car d'une part les transitions d'un état à un autre se font selon certaines probabilités et d'autre part l'émission d'un vecteur acoustique obéit à une probabilité qui dépend de l'état.

Ce modèle s'identifie tout naturellement à une chaîne de Markov qui est un processus causal.

En tout état de cause les variables de sortie de ce modèle sont supposées observables, les variables d'état ne le sont pas. Le problème sera celui de déterminer à chaque instant l'état le plus probable compte tenu d'une séquence d'observations donnée et de la connaissance supposée des paramètres du modèle.

Ce problème et le modèle qui lui est associé, justifient l'utilisation courante de l'expression anglo - saxonne de "Hidden Markov Modeling" (HMM).

Rappelons encore une fois que cette méthode ne sera pas développée. La deuxième méthode stochastique est celle que nous développerons ici. Cette nouvelle approche est appelée Modélisation des phonèmes en segments stochastiques, celle-ci est introduite pour modéliser les phonèmes à durées variables. Le phonème X est observé comme étant une séquence de vecteurs acoustiques. La durée de ce phonème X est variable selon le phonème et les cordes vocales du locuteur.

*W. H. ...*  
*...*

La modélisation des phonèmes en segments stochastiques consiste en 1)° une normalisation temporelle du segment X de longueur variable en un segment Y de longueur fixée, et en 2)° une densité multidimensionnelle des paramètres du segment normalisé Y; cette densité étant une gaussienne.

Le modèle segmental représente la structure spectrale/temporelle à travers tout le phonème. Ce modèle permet aussi l'incorporation dans Y des "caractéristiques acoustiques - phonétiques" contenues dans le segment phonétique X, ceci en plus des habituelles caractéristiques spectrales qui sont utilisées dans la méthode des chaînes de Markov cachées et la DTW.

Dans le travail qui reste, nous exposerons quelques définitions sur la segmentation de la parole, le modèle stochastique du segment phonétique, l'algorithme de reconnaissance (qui réalise simultanément la segmentation et la reconnaissance de ces derniers) et puis c'est la phase d'apprentissage dans laquelle, à partir de la parole segmentée en phonèmes donnée, nous estimerons les paramètres à savoir les moyennes et les covariances correspondants aux différents modèles phonétiques qui vont constituer le dictionnaire, puis vient la phase de ré-estimation de ces modèles en utilisant un algorithme de segmentation automatique et finalement nous présenterons un algorithme itératif qui va converger localement vers une collection optimale de modèles phonétiques.

### III.2 Segmentation de la parole

La segmentation est un des problèmes les plus difficiles à résoudre. La situation idéale serait celle où chaque segment correspondrait à un phonème. Différentes méthodes existent basées sur les courbes de variations d'énergie, ou de variabilité du signal.

Dans ce paragraphe, nous nous sommes bornés à citer brièvement les différentes approches pour la segmentation de la parole en phonèmes car dans l'algorithme de reconnaissance qu'on verra plus loin, la segmentation et la reconnaissance ont lieu simultanément.

Les différentes approches pour segmenter la parole en segments phonétiques sont :

a - Segmentation en segments d'états stables :

Cette approche repose sur le fait qu'il est possible d'identifier des segments d'état stable du signal de parole qui sont habituellement décrits en terme de changement de pression de l'air comme une fonction du temps.

b - Segmentation avec l'aide d'un modèle de référence :

Dans cette approche, le mot à segmenter est partagé en états spectraux. La plupart des phonèmes correspondent à un état spectral. A titre d'exemple, le mot "nenaanene" (mot sans signification) contient 10 états spectraux:

Le silence avant le mot /#/ , les phonèmes /n/ , /e/ , /n/ , /a/ , /n/ , /e/ , /n/ , /e/ , et le silence après le mot /#/ .

### III.3/ Modélisation des phonèmes en segments stochastiques

Cette méthode consiste en une transformation : la normalisation temporelle des segments phonétiques; et une modélisation de ces derniers par une densité gaussienne multidimensionnelle après que ces segments aient subi la normalisation temporelle.

#### - Motivations

Les principaux avantages qu'offre cette méthode sont :

1 - L'incorporation dans le segment phonétique normalisé Y, des "caractéristiques acoustiques - phonétiques" contenues dans le segment X avant normalisation (ceci en plus des habituelles caractéristiques spectrales qui sont utilisées dans la méthode des chaînes de markov cachées et la DTW).

2 - Le modèle segmental est une représentation du phonème, donc en considérant la parole à un niveau segmental, nous pouvons mieux "capturer" la structure spectrale / temporelle à travers toute la durée du phonème ceci en comparaison à ce que peut permettre une considération trame par trame.

#### III.3.1/ - Normalisation temporelle :

Soit X une suite de vecteurs acoustiques; cette séquence étant un segment de parole de longueur variable et qui représente un phonème, lequel phonème peut s'écrire comme suit :

$$X = [x_1, x_2, \dots, x_L]$$

avec  $x_i$  : Vecteur acoustique à p composantes  
L : longueur du segment phonétique

Donc X peut être considéré comme une matrice  $p \times l$  dont le nombre de colonnes est variables suivant le phonème et suivant la variabilité des cordes vocales des locuteurs.

X étant donné, on peut trouver sa représentation en longueur fixe qu'on notera Y.

Y est un segment phonétique normalisé et de longueur fixe c'est-à-dire la même longueur que tous les segments phonétiques ayant subi la normalisation temporelle et ceci quelle que soit la variabilité en longueur des phonèmes avant normalisation temporelle.

Y peut s'écrire comme suit :

$$Y = XT_L = [y_1, y_2, \dots, y_m]$$

Avec  $T_L$  : Matrice  $L \times m$  et qui permet d'effectuer la normalisation temporelle

$m$  : Longueur du segment phonétique normalisé et qui est la même pour tous les segments ayant subi la normalisation

$y_j$  : Vecteur acoustique à  $p$  composantes après normalisation temporelle (refenétrage)

Donc de même, Y peut être considéré comme une matrice  $p \times m$  dont le nombre de colonnes  $m$  est fixé ( quel que soit le nombre de colonnes de la matrice acoustique X ).

#### Définition

La normalisation temporelle consiste à choisir  $m$  instants équidistants, instants auxquels la trajectoire du segment X est refenétrée.

### La Matrice $T_L$ et ses propriétés :

Comme il a déjà été indiqué, la matrice  $T_L$  est de dimension  $L \times m$ , quant à la structure de celle-ci elle peut être trouvée comme suit :

Si on désigne par  $t_{rs}$  un élément de la matrice  $T_L$  correspondant à la  $r^{\text{ème}}$  ligne et à la  $s^{\text{ème}}$  colonne,  $t_{rs}$  peut s'exprimer comme suit :

$$t_{rs} = 1 - \alpha \quad \text{pour } r = [\beta] + 1 \quad ; \quad s = 1, 2, \dots, L \\ = \alpha \quad \text{pour } r = [\beta] + 2$$

$$\text{avec } \alpha = \beta - [\beta]; \quad \beta = \frac{(s-1)(l-1)}{(m-1)}$$

et  $[\beta]$  désigne la valeur entière de  $\beta$  n'excédant pas  $\beta$ .

Quant à l'avantage de cette transformation, il reside dans le fait que celle-ci réalise une compression de données efficace, et permet d'utiliser facilement des algorithmes de reconnaissance et de segmentation automatiques.

### III.3.2 / - le modèle probabiliste

Comme il a été déjà mentionné, le modèle segmental est une gaussienne multidimensionnelle basée sur le segment phonétique normalisé  $Y$ .

Si on désigne par  $\alpha$  un modèle phonétique contenu dans le dictionnaire du système, la densité  $P(Y/\alpha)$  ou bien d'une manière équivalente  $\ln[P(Y/\alpha)]$ , peut être interprétée comme étant le score résultant de la comparaison des caractéristiques des phonèmes à savoir le phonème normalisé  $Y$  et le modèle phonétique  $\alpha$  (qui est lui même normalisé).

Le logarithme de la probabilité conditionnelle correspondant au segment phonétique Y et au modèle phonétique  $\alpha$  peut s'exprimer d'une manière sommaire comme suit :

$$\text{Ln} [P(Y/\alpha)] = \sum_{j=1}^m \text{Ln}[P_j(y_j/\alpha)] \dots \dots \dots \quad (\text{III-1})$$

tout en sachant bien que  $y_j$  est le  $j^{\text{eme}}$  vecteur acoustique de la suite normalisée Y et m le nombre de vecteurs acoustiques constituant ladite suite.

Dans l'expression (III.1),  $P_j(y_j/\alpha)$  désigne un modèle gaussien p-dimensionnel dans lequel est mis en jeu le  $j^{\text{eme}}$  vecteur acoustique  $y_j$ , par suite on peut poser :

$P_j(y_j/\alpha) \sim N(\mu_j, C_j)$  dans laquelle N désigne la loi normale de moyenne  $\mu$  et de covariance C;  $\mu_j$  et  $C_j$  sont respectivement le  $j^{\text{eme}}$  vecteur acoustique et sa matrice covariance, lequel vecteur fait partie du modèle phonétique  $\alpha$ . Le vecteur  $\mu_j$  est mis en comparaison avec le vecteur  $y_j$  appartenant à la suite normalisée Y.

#### III.4/ - Algorithmes de reconnaissance

Dans ce paragraphe, nous décrivons l'algorithme de reconnaissance en utilisant les modèles segmentaux gaussiens. Le but de cet algorithme est de maximiser le taux de reconnaissance, lequel taux désigne la probabilité pour que le phonème reconnu  $\alpha$  soit identique au modèle phonétique  $\alpha$ . Pour accomplir ceci, nous choisissons une suite de phonèmes qui maximise la probabilité des segments phonétiques normalisés: c'est la probabilité maximum à postériori (PMA) de la suite de segments.

Plus précisément, les segments phonétiques issus de l'analyse sont premièrement transformés en segments phonétiques normalisés, puis à l'aide du calcul des probabilités de ces segments et selon le dictionnaire des modèles phonétiques, on peut trouver la probabilité maximum à postériori de la suite de segments phonétiques à reconnaître.

Avant de détailler, nous signalons que deux cas sont envisagés ici, le premier est celui où la segmentation en phonèmes est supposée réalisée préalablement, ceci uniquement dans le but de clarifier les idées. Le deuxième cas, et c'est le plus important est celui où la segmentation est supposée non réalisée, et dans lequel on généralisera en établissant un algorithme qui permettra de réaliser simultanément la segmentation en phonèmes et la reconnaissance de ces derniers: on parle de la reconnaissance automatique de la parole.

1° cas : la segmentation est supposée réalisée préalablement

Commençons par le cas le plus simple celui où la reconnaissance ne porte que sur un seul phonème. L'algorithme de reconnaissance est tout simplement la règle de la probabilité maximale à postériori dans laquelle est mis en jeu le segment normalisé Y à reconnaître :

$$\hat{\alpha} = \arg \max_{\alpha} p(\alpha/Y) \dots \dots \dots \quad (\text{III.2})$$

ou bien, d'une manière équivalente, en utilisant la règle Bayésienne, on a :

$$\hat{\alpha} = \arg \max_{\alpha} \ln [P(Y/\alpha) \cdot P(\alpha)] \quad (\text{III.2})$$

avec  $\hat{\alpha}$  : phonème reconnu

$p(\alpha)$ : probabilité à priori du modèle phonétique occurrent appartenant au dictionnaire .



Quant à l'expression  $\ln[P(Y/\alpha) \cdot P(\alpha)]$ , elle peut s'exprimer comme suit :

$$\ln[P(Y/\alpha) \cdot P(\alpha)] = \ln[P(Y, \alpha)]$$

avec :

$$\ln[P(Y, \alpha)] = -\frac{mp}{2} \ln(2\pi) - \frac{1}{2} \sum_{j=1}^m [(y_j - \mu_j(\alpha))^T \cdot (C_j(\alpha))^{-1} \cdot (y_j - \mu_j(\alpha)) + \ln|C_j(\alpha)|] \quad (\text{III.4})$$

Expression dans laquelle  $\mu_j(\alpha)$  et  $C_j(\alpha)$  sont respectivement le  $j^{\text{eme}}$  vecteur acoustique et sa matrice covariance correspondante, lequel vecteur fait partie du modèle phonétique  $\alpha$  occurrent, la notation  $|C_j(\alpha)|$  désigne le déterminant de la matrice covariance  $C_j$ .

Maximiser  $\ln[P(Y, \alpha)]$  revient à minimiser la fonction suivante:

$$D(Y, \alpha) = \sum_{j=1}^m [(y_j - \mu_j(\alpha))^T \cdot (C_j(\alpha))^{-1} \cdot (y_j - \mu_j(\alpha)) + \ln|C_j(\alpha)|] \quad (\text{III.5})$$

d'où il vient :

$$\hat{\alpha} = \arg \min_{\alpha} [D(Y, \alpha)] \quad (\text{III.6})$$

en définitive on a :

$$\hat{\alpha} = \arg \min_{\alpha} \left[ \sum_{j=1}^m [(y_j - \mu_j(\alpha))^T \cdot (C_j(\alpha))^{-1} \cdot (y_j - \mu_j(\alpha)) + \ln|C_j(\alpha)|] \right] \quad (\text{III.7})$$

Avec  $\min_{\alpha}$  désignant ici la recherche du minimum en comparant  $Y$

avec tous les modèles phonétiques ocurents contenus dans le dictionnaire.

Remarque : Les modèles phonétiques sont eux-mêmes normalisés c'est-à-dire ils sont constitués du même nombre de vecteurs acoustiques que les segments normalisés  $Y$  à reconnaître.

Maintenant supposons qu'on a à reconnaître une suite de segments phonétiques, laquelle suite sera exprimée comme suit :

$$\underline{X} = \{X_i\}_{i=1,n} \quad n : \text{Nombre de segments phonétiques}$$

Après avoir subi la normalisation temporelle, on obtient une suite de segments phonétiques normalisés, qu'on peut exprimer de la manière suivante:

$$\underline{Y} = \{Y_i\}_{i=1,n}$$

En supposant l'indépendance des segments  $X_i$  (donc des segments  $Y_i$ ), l'algorithme de reconnaissance sera celui où chaque segment  $Y_i$  subira la comparaison avec tous les modèles phonétiques occurrents du dictionnaire pour en suite trouver la probabilité maximale à postériori correspondant au phonème reconnu  $\alpha_i$ . Alors la suite de phonèmes reconnue est donnée par :

$$\underline{\hat{\alpha}} = \arg \max_{\underline{\alpha}} \{ \ln [ P(\underline{Y}/\underline{\alpha}) P(\underline{\alpha}) ] \} \quad (\text{III.8})$$

avec  $\underline{\hat{\alpha}} = \{\alpha_i\}_{i=1,n}$

$P(\underline{\alpha})$  désigne les probabilités à priori des modèles phonétiques qui subissent la comparaison avec les phonèmes normalisés qui constituent la suite  $\underline{Y}$ .

$$\text{Ln}[P(\underline{Y}/\underline{\alpha})] = \sum_{i=1}^n \text{Ln}[P(Y_i/\alpha_i)P(\alpha_i)] \quad (\text{III.9})$$

L'indice  $i$  dans la notation " $\alpha_i$ " veut dire tout simplement que l'ensemble des modèles phonétiques sont comparés avec le  $i$ ème segment à savoir  $Y_i$ , ceci d'une part, d'autre part la suite de phonèmes reconnue peut s'écrire de la façon suivante :

$$\hat{\underline{\alpha}} = \arg \max_{\underline{\alpha}} \left[ \sum_{i=1}^n \text{Ln}[P(Y_i/\alpha_i)P(\alpha_i)] \right] \quad (\text{III.10})$$

ou bien :

$$\hat{\underline{\alpha}} = \arg \min_{\underline{\alpha}} \left[ \sum_{i=1}^n D(Y_i, \alpha_i) \right] \quad (\text{III.11})$$

où  $D(Y_i, \alpha_i)$  est la fonction définie par (III.5)

2° cas: la segmentation n'est pas réalisée préalablement

Soit  $\underline{X}$  la suite des phonèmes à reconnaître qui peut s'écrire comme auparavant:  $\underline{X} = \{X_i\}_{i=1,n}$

après normalisation, on obtient une suite normalisée  $\underline{Y}$  :

$$\underline{Y} = \{Y_i\}_{i=1,n}$$

Si on désigne par  $\underline{x}$  la suite de tous les vecteurs acoustiques constituant  $\underline{X}$ , alors on peut écrire:

$$\underline{x} = \{x_1\}_{1=1,N} ; \quad N: \text{Nombre total des vecteurs acoustiques constituant } \underline{X}$$

la segmentation correspondant à la séquence  $\underline{x}$ , peut être représentée par :  $\underline{s} = \{ s(i) \}_{i=1,n}$  dans laquelle  $s(i)$

désigne l'indice du dernier vecteur acoustique contenu dans le  $i$ ème segment phonétique  $X_i$  lequel, (en utilisant cette notation) peut s'exprimer de la manière suivante :

$$X_i = \{ x_{s(i-1)+1}, \dots, x_{s(i)} \}$$

Avec la supposition que  $s(0)=0$

$X_i$  est un segment phonétique de longueur  $L(i) = s(i) - s(i-1)$  et bien sûr on a :  $Y_i = X_i^T L(i)$ .

Rappelons que ce qui vient d'être exposé n'est qu'une approche purement mathématique de la segmentation et que cela va servir dans ce qui va suivre.

En reconnaissance automatique de la parole, on détermine la segmentation du signal parole d'entrée et on reconnaît la suite de phonèmes qui en correspond.

L'approche Bayésienne pour trouver la vraisemblance de la suite de phonèmes est donnée par :

$$L_B(\underline{\alpha}) = \sum_{\underline{s}} [P(\underline{Y}(\underline{s})/\underline{\alpha}) P(\underline{\alpha})]$$

(III.12)

expression dans laquelle la notation  $\underline{Y}(s)$  exprime tout simplement le fait que la suite  $\underline{Y}$  dépend de la segmentation. Comme cette approche est très coûteuse en temps de calcul, on lui substitue la méthode de détection par le maximum de vraisemblance lorsque certains paramètres du signal sont inconnus. Nous considérons que la segmentation est un paramètre inconnu.

Pour chaque suite hypothétique de phonèmes, on détermine la plus vraisemblable segmentation en maximisant la vraisemblance:

$$\underline{L}(\underline{\alpha}) = \max_{\underline{s}} \{ \text{Ln} [ P(\underline{Y}(\underline{s}) / \underline{\alpha}) \cdot P(\underline{\alpha}) ] \} \quad (\text{III-13})$$

alors la suite de phonèmes reconnue est :

$$\hat{\underline{\alpha}} = \arg \max_{\underline{\alpha}} (\underline{L}(\underline{\alpha})) \quad (\text{III-14})$$

l'algorithme peut se résumer ainsi:

1- hypothétiser toutes les segmentations possibles:

$$\underline{s} = \{ s(i) \} \quad i=1, n$$

2- normaliser toutes les segmentations:

$$X_i = [ x_{s(i-1)+1}, \dots, x_{s(i)} ]$$

$$L(i) = s(i) - s(i-1)$$

$$Y_i = X_i \cdot T_{L(i)}$$

3- trouver la meilleure suite de phonèmes et faire correspondre une vraisemblance pour chaque segmentation hypothétique:

$$\underline{L}(\underline{\alpha}) = \max_{\underline{s}} \{ \text{Ln} [ P(\underline{Y}(\underline{s}) / \underline{\alpha}) \cdot P(\underline{\alpha}) ] \}$$

4- Parmi toutes les vraisemblances trouvées, choisir celle qui est maximale, laquelle vraisemblance correspond à la suite de phonèmes reconnue :

$$\hat{\underline{\alpha}} = \arg \max_{\underline{\alpha}} \mathcal{J}(\underline{\alpha})$$

Pour des raisons pratiques, (III-13) n'est pas utilisée alors on lui substitue une expression dans laquelle figure une somme de probabilités logarithmiques pondérées. <sup>Cette</sup> ~~laquelle~~ somme exprime le score résultant de la comparaison des caractéristiques des phonèmes (les phonèmes à reconnaître et les modèles phonétiques) dans ce cas le score pour la segmentation à maximiser à partir de l'ensemble des segmentations permises sera donné par :

$$J(\underline{Y}/\underline{\alpha}) = \sum_{i=1}^n \{ \text{Ln}[P(Y_i/\alpha_i)] \cdot L(i) + \text{Ln}[P(\alpha_i)] + C \} \quad (\text{III-15})$$

où n désigne le nombre de segments.

Il est évident que  $J(\underline{Y}/\underline{\alpha})$  remplace  $\text{Ln}[P(\underline{Y}(\underline{s})/\underline{\alpha}) \cdot P(\underline{\alpha})]$  dans (III-13) alors celle-ci s'écrira :

$$\mathcal{L}(\underline{\alpha}) = \max_{\underline{s}} \{ \sum \{ \text{Ln}[P(Y_i/\alpha_i)] \cdot L(i) + \text{Ln}[P(\alpha_i)] + C \} \} \quad (\text{III-16})$$

par conséquent la suite de phonèmes reconnue s'écrira :

$$\hat{\underline{\alpha}} = \arg \max_{\underline{\alpha}} \{ \max_{\underline{s}} \{ \sum \{ \text{Ln}[P(Y_i/\alpha_i)] \cdot L(i) + \text{Ln}[P(\alpha_i)] + C \} \} \} \quad (\text{III-17})$$

commentaire :

Le paramètre C désigne le coût par segment. Il est utilisé pour ajuster le taux d'insertion des phonèmes dans le système. ce

coût correspond à la probabilité logarithmique du taux de phonèmes ( nombre de phonèmes/emission )

D'autre part la durée  $L(i)$  (ou longueur) du  $i^{\text{ème}}$  phonème (à reconnaître) est introduite uniquement pour s'assurer qu'un segment phonétique  $X_i$  de longueur  $L(i)$  contribue au score dans la même proportion que sa longueur .

#### Autre Algorithme:

En plus de l'algorithme déjà décrit ,il existe une autre solution qui utilise un algorithme de programmation dynamique. Plus précisément , à chaque temps  $t$  , on calcule le score pour la meilleure suite de phonèmes reconnue qui prend fin à l'instant  $t$  :

$$(J^*)_t = \max_{\tau, \alpha} \{ (J^*)_{\tau} + \text{Ln}[P(Y(\tau, t)/\alpha)] \cdot (t - \tau) + \text{Ln}[P(\alpha)] + C \} \quad (\text{III-18})$$

$Y(\tau, t) = [ X_{\tau+1}, \dots, X_t ]$   $T_{t-\tau}$  et  $(J^*)_t$  est le score de la meilleure suite de phonèmes reconnue . Le dernier phonème  $\alpha^*$  correspondant à  $(J^*)_t$ , ainsi que l'instant  $\tau^*$  où prend fin le phonème qui le précède sont notés pour pouvoir déterminer le ensuite la séquence de phonèmes reconnue dans sa totalité . La solution à la fin de la phrase, à l'instant  $t_f$ , est donnée par  $(J^*)_{t_f}$ . Cet algorithme est très coûteux en temps de calcul que ne l'est l'algorithme précédent. Toutefois la complexité de la recherche est proportionnelle au nombre de modèles phonétiques contenus dans le dictionnaire .

#### III-5/ PHASE D'APPRENTISSAGE

Dans ce paragraphe, nous abordons l'entraînement des modèles phonétiques . Nous commençons par décrire une procédure pour l'estimation statistique des modèles phonétiques à partir

d'une séquence de parole préalablement segmentée en phonèmes. Ensuite, nous décrirons un algorithme qui va réaliser la segmentation automatique laquelle segmentation va servir de ré-estimation des modèles phonétiques résultant de l'estimation tout en ayant leur(s) transcriptions phonétiques, puis à partir de ces deux étapes on établira un algorithme itératif qui réalisera la segmentation automatique à partir de la parole continue.

### III-5-1/ ESTIMATION DES PARAMETRES

A partir d'une séquence de parole segmentée en phonèmes il est possible d'estimer les modèles phonétiques à partir des statistiques des segments phonétiques normalisés. Les probabilités à priori des phonèmes sont estimées à partir de leur fréquences relatives dans la séquence de la parole d'entraînement.

Si on suppose que pour un segment normalisé tous les vecteurs acoustiques qui le constituent sont indépendants, alors un modèle phonétique est donné par les  $m$  vecteurs acoustiques modèles indépendants.

Par conséquent, il est simplement nécessaire d'estimer le modèle gaussien pour chaque vecteur acoustique et pour chaque modèle phonétique.

Cette estimation porte sur les statistiques des vecteurs acoustiques constituant les segments phonétiques normalisés donnés.

La densité du  $j^{\text{ième}}$  vecteur acoustique du phonème  $\alpha$  est désignée par  $N(\mu_j, C_j)$ . Cette densité est déterminée par les méthodes d'estimation classiques pour la moyenne et la covariance (voir annexe) ; ceci à partir de l'ensemble des vecteurs acoustiques qui correspondent au  $j^{\text{ième}}$  vecteur du phonème .



### III-5-2/ SEGMENTATION AUTOMATIQUE :

A partir de la collection des modèles gaussiens et leurs transcriptions phonétiques on peut déterminer la segmentation d'une séquence de parole de transcription phonétique  $\alpha$  connue. Cette phase qui utilise un algorithme très simple permet de ré-estimer les modèles phonétiques obtenus à l'étape de l'estimation, puis on établira un algorithme itératif qui permettra d'entraîner automatiquement les modèles phonétiques à partir de la parole continue .

L'algorithme est simplement une recherche parmi toutes les segmentations possibles et ensuite choisir celle qui donne la transcription  $\alpha$  avec un maximum de probabilité .

Cette segmentation s'écrit :

$$\underline{s} = \arg \max_{\underline{s}} \mathcal{L}(\underline{s}) \quad (\text{III-19})$$

où la vraisemblance de segmentation est donné par :

$$\mathcal{L}(\underline{s}) = \text{Ln}[P(\underline{Y}(\underline{s})/\underline{\alpha}) \cdot P(\underline{\alpha})] \quad (\text{III-20})$$

$$= \sum_{i=1}^n \{ \text{Ln}[P(Y_i/\alpha_i)] \cdot L(i) + \text{Ln}[P(\alpha_i)] + C \} \quad (\text{III-21})$$

alors, la segmentation est donnée par :

$$\hat{\underline{s}} = \arg \max_{\underline{s}} \left\{ \sum_{i=1}^n \{ \text{Ln}[P(Y_i/\alpha_i)] \cdot L(i) + \text{Ln}[P(\alpha_i)] + C \} \right\} \quad (\text{III-22})$$

Comme en phase de reconnaissance, le score  $f(\underline{Y}/\underline{\alpha})$  (défini en III-15) remplace l'expression  $\ln [P(\underline{Y}(\underline{s})/\underline{\alpha}) P(\underline{\alpha})]$ , ceci dans le but d'introduire le facteur de coût et la durée (ou longueur) de chaque phonème dans la séquence de parole .

L'introduction de ces deux facteurs va améliorer la probabilité maximum correspondant à la segmentation  $\hat{\underline{s}}$  .

#### Algorithme d'entraînement Itératif:

En utilisant les deux étapes déjà décrites , l'estimation des paramètres et la segmentation automatique, nous pouvons définir un algorithme itératif pour entraîner automatiquement les modèles segmentaux à partir de la parole continue.

si on se donne:

\* La transcription phonétique de la séquence d'entraînement

\* Les modèles gaussiens initiaux pour tous les phonèmes:

{  $P_0(Y/\alpha)$  }

\*  $t=0$

On itère:

1) déterminer la segmentation à probabilité maximale  $\hat{\underline{s}}_t$  correspondante à la séquence d'entraînement pour la transcription donnée et les densités de probabilité courantes {  $P_t(Y/\alpha)$  }

2) déterminer l'estimation par le maximum de vraisemblance pour les densités {  $P_{t+1}(Y/\alpha)$  } en utilisant  $\hat{\underline{s}}_t$ .

3)  $t \leq t+1$  et aller à l'étape 1)

Il apparaît qu'à chaque étape de l'algorithme la vraisemblance de la segmentation  $\mathcal{L}(\hat{\underline{s}}_t)$  croît pour la séquence de phonèmes (d'entraînement) donnée.

Dans l'étape 1), la nouvelle segmentation est la plus vraisemblable pour les modèles phonétiques courants,

qui devrait être, au moins, aussi probable que la segmentation qui précède. Dans l'étape 2), les nouvelles densités font croître la probabilité de la segmentation courante ceci par définition de l'estimation par le maximum de vraisemblance.

D'autre part, le nombre de segmentations est fini car la longueur de la séquence d'entraînement est aussi finie, donc  $\mathcal{L}(\hat{s}_t)$  est borné par la vraisemblance de la meilleure segmentation. Par conséquent, la vraisemblance de la segmentation  $\mathcal{L}(\hat{s}_t)$  converge vers un optimum local. La séquence des segmentations converge vers un optimum local quand la vraisemblance de la segmentation de la séquence l'est du fait que le nombre de segmentations possibles est fini et la recherche de la vraisemblance maximale dans l'étape 1) est imposée.

Ce qui conduit vers une collection optimale de modèles phonétiques.

## **ANNEXES**

## A. NOTIONS DE LA THEORIE DES PROBABILITES

A - Rappel des définitions et des théorèmes élémentaires:

a - Soient tous les évènements mutuellement exclusifs d'un espace d'échantillonnage alors:

$$\sum_i P(A_i) = 1$$

b - Probabilité conditionnelle: on considère deux évènements A et B. La probabilité conditionnelle  $P(B/A)$  est la probabilité de réalisation de B sachant que A s'est réalisé.

c - probabilité composée: C'est la probabilité de réalisation simultanée des évènements A et B:

$$P(A,B) = P(A) \cdot P(B/A) = P(B) \cdot P(A/B)$$

Pour plusieurs événements simultanées:

$$P(A_1, B_2, \dots, A_n) = P(A_1) \cdot P(A_2/A_1) \cdot \dots \cdot P(A_n/A_1, \dots, A_{n-1})$$

d - Evénements indépendants: deux événements sont dits indépendants si:

$$P(A/B) = P(A)$$

Théorème: Soient deux événements indépendants A et B alors:

$$P(A, B) = P(A) \cdot P(B)$$

e - Théorème des probabilités totales: on considère les événements  $H_i$  mutuellement exclusifs dont l'union est l'espace fondamental ... si A est un événement quelconque alors:

$$P(A) = \sum_i P(H_i) \cdot P(A/H_i)$$

d - Théorème de Bayes

on prend les mêmes conditions que le théorème précédent:

$$P(H_i/A) = P(H_i) \cdot P(A/H_i) / P(A)$$

A - 2 Variable Aléatoire, Loi de distribution:

a - Une variable X est dite aléatoire si elle prend ses valeurs selon une probabilité P(X).

Une variable aléatoire X peut être discrète ou continue .

Il est commode d'introduire la fonction de probabilité ou fonction de répartition P(X=x).

Dans le cas d'une variable aléatoire discrète on a bien sûr:

$$\sum_k P(X=x_k) = 1$$

( $x_k$ : chacune des valeurs que peut prendre la variable aléatoire).

Pour une variable aléatoire continue on préfère définir une fonction densité de probabilité  $P(x)$  telle que:

$$P(x) \geq 0$$

$$\int_{-\infty}^{+\infty} P(x) dx = 1$$

b - Loi normal ou loi de gauss:  $P(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-u)^2}{2\sigma^2}\right)$

$u$ : moyenne ou espérance de  $X$ ;  $\sigma^2$ : variance loi normal multidimensionnelle pour des variables indépendants:

$$P(x_1, \dots, x_n) = \frac{1}{\pi^{n/2} \prod_{i=1}^n \sigma_i} \exp\left(-\sum_{i=1}^n \frac{(x_i - u_i)^2}{2\sigma_i^2}\right)$$

## B - STATISTIQUES

a - Estimation des moyennes:

Soient le vecteur aléatoire suivant:  $X = (X_1, X_2, \dots, X_k)$

où  $X_i$  est une variable aléatoire.

Soit  $N$  le nombre de valeurs que prend ce vecteur aléatoire:

$$x_1 = [ (x_1)_1, (x_2)_1, (x_3)_1, \dots, (x_k)_1 ]$$

$$x_2 = [ (x_1)_2, (x_2)_2, (x_3)_2, \dots, (x_k)_2 ]$$

⋮

$$x_N = [ (x_1)_N, (x_2)_N, (x_3)_N, \dots, (x_k)_N ]$$

alors la moyenne correspondante à la variable aléatoire  $X_i$ :

$$\mu_{X_i} = \frac{\sum f(x_j)_i}{N} \quad \text{avec } (x_j)_i = \text{valeur de la variable } x_j \text{ correspondant à la } i\text{ème réalisation}$$

on aura alors le vecteur moyenne:

$$\mu = (\mu_{X_1}, \mu_{X_2}, \mu_{X_3}, \dots, \mu_{X_i}, \dots, \mu_{X_k})$$

b - Estimation de la matrice covariance:

Soit  $X$  le vecteur aléatoire définie précédemment, alors sa matrice covariance s'écrira:

$$C_X = \begin{bmatrix} \sigma^2_{X_1} & c_{X_1X_2} & \dots & c_{X_1X_k} \\ c_{X_2X_1} & \sigma^2_{X_2} & & \\ \cdot & & \cdot & \\ \cdot & & & \sigma^2_{X_k} \\ c_{X_kX_1} & & & \end{bmatrix}$$



avec  $C_{x_i x_j} = \frac{1}{N} \sum_{l=1}^N [(x_i)_l - \mu_{x_i}] [(x_j)_l - \mu_{x_j}] / N$

et  $C_{x_i x_i} = \sigma^2_{x_i} = \frac{1}{N} \sum_{l=1}^N [(x_i)_l - \mu_{x_i}]^2 / N$

Où  $\mu_{x_i}$ ,  $\mu_{x_j}$  désignent les moyennes correspondantes aux variables adjacentes  $x_i$ ,  $x_j$ .

### C - RESULTATS DE LA SIMULATION

#### a - Analyse:

La méthode d'analyse utilisée est la méthode par prédiction linéaire. Le signal parole simulé renferme les cinq voyelles: a, o, i, e, u, chacune de ces voyelles est générée par une somme de six sinusoïdes où sont exprimés l'amplitude et les deux formants de la voyelle.

La durée de la voyelle intervient dans les phases de l'échantillonnage et du fenêtrage.

Toutes ces opérations sont effectuées dans les conditions suivantes:

- Fréquence d'échantillonnage:  $f_e = 12.8 \text{ KHZ}$

- Nombre d'échantillons de la fenêtre d'analyse :

$$N = 256$$

- Ordre de prédiction linéaire :  $P = 12$

Caractéristiques des voyelles générées:

pour le phonème a:

amplitude: 80 mv

premier formant: 750 Hz

deuxième formant: 1350 Hz

durée : 0.8 ms

Pour le phonème o:

amplitude : 80 mv  
premier formant: 375 Hz  
deuxième formant: 750 Hz  
durée: 0.8 ms

Pour le phonème i:

amplitude: 80 mv  
premier formant: 250 Hz  
deuxième formant: 2500 Hz  
durée: 0.8 ms

Pour le phonème e:

amplitude: 80 mv  
premier formant: 400 Hz  
deuxième formant: 2200 Hz  
durée: 0.8 ms

Pour le phonème u:

amplitude: 80 mv  
premier formant: 250 Hz  
deuxième formant: 600 Hz  
durée: 0.8 ms

b - Reconnaissance:

Dans cette phase, on a d'abord créé un dictionnaire qui contient la collection des cinq modèles phonétiques (voyelles) où chaque modèle est créé en estimant les moyennes et les covariances des différents vecteurs acoustiques correspondants aux différentes versions du phonème. Ces versions sont obtenues en perturbant légèrement les caractéristiques des phonèmes; puis vient l'étape de ré-estimation dans laquelle on cherche à reconnaître chaque phonème individuellement tout en perturbant légèrement ses caractéristiques, nous avons répété plusieurs

fois cette opération de reconnaissance jusqu'à ce que la reconnaissance reste parfaite malgré les légères perturbations qu'on subit les caractéristiques à chaque répétition. Toutefois les perturbations ne devront pas dépasser une certaine limite pour rester toujours dans le cadre d'un système phonatoire humain. Une fois l'opération de ré-estimation s'est soldée par de bons résultats, on a pris comme optimaux les modèles estimés statistiquement (par les moyennes et les covariances).

Nous précisons que chaque modèle optimal est affectée d'une étiquette qui va servir comme moyen de décodage dans la reconnaissance simulée.

Concernant la reconnaissance simulée proprement dite, les mots à reconnaître sont des suites de phonèmes (ici des voyelles),

Nous avons utilisé les étiquettes suivantes:

phonème	étiquette
a	"a"
o	"o"
i	"i"
e	"e"
u	"u"

nous avons effectué les tests suivants où chaque mot est représenté par une suite de cinq phonèmes (taux d'insertion)

TEST1

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'aoieu'

Donner les amplitudes des phonèmes successifs:

80 80 80 80 80 - du CM

Donner les premiers formants des phonèmes successifs:

750 375 250 400 250

Donner les deuxièmes formants des phonèmes successifs:

1350 750 2500 2200 600

Donner les durées des phonèmes successifs:

0.8 0.8 0.8 0.8 0.8

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'aoieu'

TEST2

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'ueioa'

Donner les amplitudes des phonèmes successifs:

80.5 80.9 80.7 80.6 81

Donner les premiers formants des phonèmes successifs:

255 405 255 380 755

Donner les deuxièmes formants des phonèmes successifs:

605 2205 2505 755 1355

Donner les durées des phonèmes successifs:

0.81 0.82 0.83 0.84 0.85 ms

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'ueioa'

TEST3

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'ioaeu'

Donner les amplitudes des phonèmes successifs:

79.5 78.5 82 81 83

Donner les premiers formants des phonèmes successifs:

260 385 760 410 260

Donner les deuxièmes formants des phonèmes successifs:

2510 760 1360 2210 610

Donner les durées des phonèmes successifs:

0.75 0.78 0.74 0.82 0.82

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'ioaeu'

TEST4

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'euaoi'

Donner les amplitudes des phonèmes successifs:

75 76 77 78 80

Donner les premiers formants des phonèmes successifs:

390 240 760 365 240

Donner les deuxièmes formants des phonèmes successifs:

2190 590 1340 760 2490

Donner les durées des phonèmes successifs:

0.75 0.76 0.77 0.78 0.79

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'euaoi'

## TESTS

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'oieua'

Donner les amplitudes des phonèmes successifs:

75 76 78 77 81

Donner les premiers formants des phonèmes successifs:

370 245 395 245 745

Donner les deuxièmes formants des phonèmes successifs:

745 2495 2195 595 1345

Donner les durées des phonèmes successifs:

0.75 0.78 0.76 0.74 0.73

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'oieua'



## TEST6

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'ieuao'

Donner les amplitudes des phonèmes successifs:

75 76 80 82 83

Donner les premières formants des phonèmes successifs:

253 405 255 755 320

Donner les deuxièmes formants des phonèmes successifs:

2505 2206 608 1356 755

Donner les durées des phonèmes successifs:

0.81 0.32 0.79 0.75 0.78

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédictifs:

12

Le mot reconnu est: 'euiaa'

TEST7

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'aiuce'

Donner les amplitudes des phonèmes successifs:

80 80.4 79.6 78.8 77.8

Donner les premiers formants des phonèmes successifs:

755 254 257 380 408

Donner les deuxièmes formants des phonèmes successifs:

1957 2508 610 759 2208

Donner les durées des phonèmes successifs:

0.78 0.79 0.75 0.76 0.81

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'aiuce'

TEST8

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'eueue'

Donner les amplitudes des phonèmes successifs:

80.2 80.3 80.4 78.25 75.66

Donner les premiers formants des phonèmes successifs:

400 250 405 255 412

Donner les deuxièmes formants des phonèmes successifs:

2205 600 2206 610 2209

Donner les durées des phonèmes successifs:

0.78 0.79 0.75 0.81 0.80

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'eueue'

TEST9

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'auaua'

Donner les amplitudes des phonèmes successifs:

80.4 80.56 78.5 79.55 82.1

Donner les premiers formants des phonèmes successifs:

750 250 755 255 746

Donner les deuxièmes formants des phonèmes successifs:

1350 605 1348 609 1346

Donner les durées des phonèmes successifs:

0.8 0.75 0.79 0.76 0.78

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'auaua'

TEST10

Donner le nombre de phonèmes du mot à reconnaître:

5

Donner la suite des phonèmes du mot:

'iiii'

Donner les amplitudes des phonèmes successifs:

70 71 72 73 74

Donner les premiers formants des phonèmes successifs:

230 220 200 205 206

Donner les deuxièmes formants des phonèmes successifs:

2500 2400 2350 2300 2250

Donner les durées des phonèmes successifs:

0.70 0.71 0.72 0.73 0.74

Donner le nombre de vecteurs acoustiques avant la normalisation:

40

Donner le nombre de vecteurs acoustiques après la normalisation:

10

Donner le nombre de coefficients prédicteurs:

12

Le mot reconnu est: 'iiie'

## Interprétation des résultats des tests:

→ A travers les résultats de ces tests, on voit bien l'efficacité du caractère statistique de la méthode car en variant les caractéristiques des phonèmes dans une plage de  $\pm 1.5\%$  des valeurs de référence, on continue quand-même à reconnaître le mot. Ce fait confère à la méthode d'être utilisée dans un système de reconnaissance multilocuteur.

→ d'autre part, nous avons constaté que le temps de reconnaissance est très important et il serait d'autant plus important que le nombre de modèles phonétiques est grand ce qui nous permet bien de faire une suggestion en préconisant d'implanter la méthode dans des microprocesseurs puissants dans le but d'augmenter la vitesse de reconnaissance qui est l'un des plus grands objectifs en reconnaissance de la parole.

19  
26

Pour finir, on peut ajouter que dans cette méthode la reconnaissance porte sur l'unité fondamentale indivisible de la parole qu'est le phonème, donc le taux d'erreurs de reconnaissance est très faible ce qui n'est pas le cas avec les méthodes globales.