



UNIVERSITÉ DE TOULOUSE III

TRAITEMENT AUTOMATIQUE DE LA PAROLE

Système de vérification de locuteurs

Auteurs :

M. Axel BELLEC
M. Max HALFORD

Enseignant :

M. Jérôme FARINAS
M. Julien PINQUIER

Table des matières

1	Présentation de la base de données	2
2	Décomposition parole/non-parole	3
3	Paramétrisation	5
4	Apprentissage	7
5	Reconnaissance	8
6	Affichage des résultats	9
7	Fonction main	11

1 Présentation de la base de données

Ecrire une fonction «lecture.m» permettant de lire (charger dans une variable) sous Matlab un fichier son et de connaître ses caractéristiques (fréquence d'échantillonnage, nombre de bits de quantification, durée).

```
1 function [signal, fs, nb_bits, duration]
   ↪ = lecture(fichier)
2     [signal, fs] = audioread(fichier);
3     nb_bits = 8;
4     duration = nb_bits * fs;
5 end
```

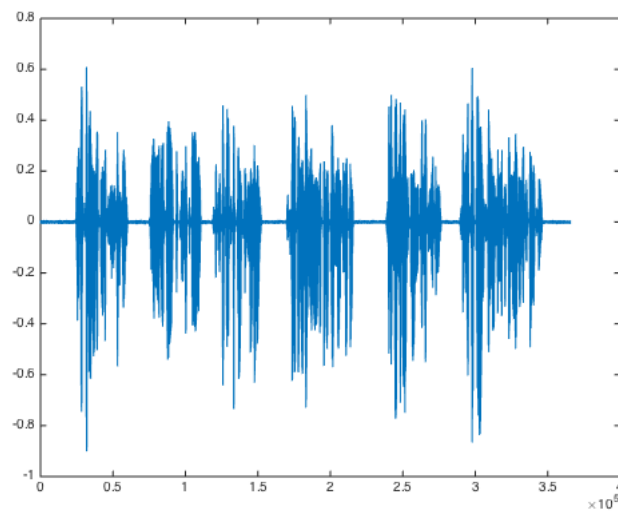


FIGURE 1 – Représentation graphique du signal

```
1 signal = lecture('APP/L1_fic1.wav');
2 plot(signal);
```

2 Décomposition parole/non-parole

Ecrire une première fonction «energie.m» permettant de calculer l'énergie à court terme d'un signal en fonction de la taille des fenêtres d'analyse.

Le signal étant représenté par un vecteur de longueur n , on va tout simplement boucler de 1 à n avec un pas de $\left\lfloor \frac{n}{p} \right\rfloor$ où p est le nombre d'intervalles voulus. Evidemment le dernier intervalle aura $n \bmod p$ valeurs si n n'est pas un multiple de p .

```
1 function [energies] = energie(signal, taille_fenetre)
2     energies = [];
3     n = length(signal);
4     pas = floor(taille_fenetre / 2);
5     for i = 1:pas:n
6         energie = sum(signal(i:i+pas-1)).^2) /
↪     taille_fenetre;
7         fprintf('from: %d | until: %d | energie:
↪     %d.\n', from, until, energie);
8         energies = [energies, energie];
9     end
10 end
```

Ecrire une autre fonction «etiquetage.m» (utilisant la fonction «energie») permettant d'étiqueter un fichier son en parole et non-parole sous Matlab.

Pour étiqueter un fichier son en zone de parole/non-parole, il suffit de parcourir le vecteur de l'énergie du signal. On regarde ensuite si pour chaque valeur d'énergie est supérieure à notre seuil. Si c'est le cas, on considère que c'est une zone de parole.

```
1 function [etiq_parole] = etiquetage(signal,
↪     taille_fenetre, threshold)
2     etiq_parole = [];
3     [energies] = energie(signal, taille_fenetre);
4     for i=1:length(energies)
5         if energies(i) >= threshold
6             etiq_parole = [etiq_parole, 1];
7         else
8             etiq_parole = [etiq_parole, 0];
9         end
10    end
11 end
```

Ecrire une fonction «etiquetage_total.m» (utilisant la fonction «etiquetage») permettant d'étiqueter en parole/non-parole l'ensemble des fichiers sonores

du répertoire «WAV/APP».

On peut itérer à travers tout les fichiers du répertoire en remarquant qu'ils sont numérotés. On peut simplement faire boucler i de 1 au nombre de locuteurs et j de 1 au nombre de fichiers par locuteurs pour ensuite interpoler i et j afin d'obtenir le chemin de chaque fichier.

```
1 function etiquetage_total(taille_fenetre, nbe_loc,
   ↪ nbe_fic, threshold)
2     basepath = 'APP';
3     output_path = 'LABELS';
4     % Loop through nbe_loc
5     for i=1:nbe_loc
6         % And then loop through nbe_fic
7         for j=1:nbe_fic
8             % Generate audio filename
9             input_file = sprintf('%s/L%d_fic%d.wav',
   ↪ basepath, i, j);
10            [signal, fs] = lecture(input_file);
11            [etiq_parole] = etiquetage(signal,
   ↪ taille_fenetre, threshold);
12            output_file = sprintf('%s/L%d_fic%d.lab',
   ↪ output_path, i, j);
13            save(output_file, 'etiq_parole', '-ASCII');
14        end
15    end
16 end
```

3 Paramétrisation

Ecrire un programme permettant d'effectuer la paramétrisation (calcul des coefficients cepstraux).

Le cepstre est le résultat de la transformée de Fourier appliquée au logarithme naturel de la transformée de Fourier du signal dont la phase est ignorée. Ici on applique tout simplement le schéma classique d'extraction indiqué dans le TP.

```
1 function [mfcc] = parametrisation(signal,
   ↪  taille_fenetre, nbe_coef)
2     mfcc = []; % Initialize matrix
3     n = length(signal);
4     pas = floor(taille_fenetre / 2);
5     for i=1:pas:n
6         cepstre_real =
   ↪  real(ifft(log(abs(fft(signal(i:(i+pas-1))) .*
   ↪  hamming(pas))))));
7         calc_mfcc = cepstre_real(2:nbe_coef);
8         mfcc = [mfcc; calc_mfcc']; % Transpose
   ↪  calc_mfcc
9     end
10 end
```

Ecrire une fonction «parametrisation_total.m» (utilisant la fonction «parametrisation») permettant de calculer les MFCC pour l'ensemble des fichiers du répertoire «WAV/APP»

De la même manière que pour la fonction `etiquetage_total`, nous allons itérer sur les locuteurs i et leurs fichiers associés j pour calculer leur MFCC. Nous enregistrons ensuite la matrice obtenue dans un fichier.

```

1 function parametrisation_total(taille_fenetre,
  ↳ nbe_loc, nbe_fic, nbe_coef)
2     basepath = 'APP';
3     output_path = 'MFCC';
4     % Loop through nbe_loc
5     for i=1:nbe_loc
6         % And then loop through nbe_fic
7         for j=1:nbe_fic
8             % Generate audio filename
9             input_file = sprintf('%s/L%d_fic%d.wav',
  ↳ basepath, i, j);
10            [signal, fs] = lecture(input_file);
11            mfcc = parametrisation(signal,
  ↳ taille_fenetre, nbe_coef);
12            output_file =
  ↳ sprintf('%s/L%d_fic%d.mfcc', output_path, i, j);
13            save(output_file, 'mfcc', '-ASCII');
14        end
15    end
16 end

```

4 Apprentissage

Ecrire la fonction «apprentissage.m» permettant d'effectuer l'apprentissage des modèles. Cette fonction prendra en entrée les vecteurs de données à traiter (fichiers «monde» et «locuteur») ainsi que le nombre de lois gaussiennes n . Elle fournira en sortie les modèles du monde et du locuteur représentés par les poids w , les moyennes m et les variances v de chacune des lois gaussiennes

Ici on applique l'algorithme mis à disposition pour déterminer les paramètres des gaussiennes que l'on va utiliser pour modéliser le “monde” et les locuteur. Les paramètres sont d'abord approximés grâce à l'algorithme des K-means. Pour avoir des paramètres plus réalistes, l'algorithme EM est utilisé en calculant la probabilité d'appartenance de chaque locuteur à chaque distribution (*Expectation*) et en re-paramétrisant les distributions (*Maximization*).

```
1 function [m, v, w, m_l, v_l, w_l]
   ↪ = apprentissage(fic_monde, fic_locuteur, n)
2     monde = load(fic_monde);
3     [m, v, w] = gaussmix(monde, 1, 10.001, n);
4     locuteur = load(fic_locuteur);
5     [m_l, v_l, w_l] = gaussmix(locuteur, 1, 10.001,
   ↪ m, v, w);
6 end
```


5 Reconnaissance

Ecrire la fonction «*tests_total.m*» permettant d'effectuer la reconnaissance des fichiers sons inconnus situés dans le répertoire «WAV/RECO». Cette fonction prendra en entrée la taille de la fenêtre d'analyse, le nombre de coefficients cepstraux ainsi que les paramètres (poids, moyennes et matrices de covariance) de chacun des modèles, composés de mélanges de lois gaussiennes. Elle fournira en sortie un vecteur (ou matrice) du taux de reconnaissance par fichier.

On va comparer les probabilités d'appartenance pour un locuteur à la distribution de la gaussienne “monde” et à celle du locuteur. On utilise le MAP pour décider quelle distribution choisir (tout simplement on garde la probabilité retournée par chaque distribution la plus élevée).

```
1 function [taux_reco] = tests_total(taille_fenetre,
   ↪ nbe_coef, m, v, w, m_l, v_l, w_l)
2     reco_path = 'RECO';
3     taux_reco = [];
4     for i=1:10
5         for j=9:10
6             input_file = sprintf('%s/L%d_fic%d.wav',
   ↪ reco_path, i, j);
7             [signal, fs] = lecture(input_file);
8             % Compute MFCC for test files
9             mfcc = parametrisation(signal,
   ↪ taille_fenetre, nbe_coef);
10            % Compute probabilities for each window
11            prob_monde = gmm1pdf(mfcc, m, v, w);
12            prob_locu = gmm1pdf(mfcc, m_l, v_l, w_l);
13            % Create a new vector with 1 if locutor
   ↪ wins, else 0
14            prob_tot = prob_monde < prob_locu;
15            % Compute the sum for locutor
16            locu_wins = sum(prob_tot);
17            % Compute the final score
18            score = locu_wins / length(prob_tot);
19            taux_reco = [taux_reco, score];
20        end
21    end
22 end
```

6 Affichage des résultats

Nous avons développé une procédure supplémentaire pour afficher la classification que nous avons effectuée en fonction de la vraie classe.

```
1 function print_truth_predict(taux_reco)
2     files = cell(20,1);
3     locutor_truth = cell(20,1);
4     basepath = 'REC0';
5     index = 1;
6     for i=1:10
7         for j=9:10
8             files{index} =
9             ↪ sprintf('%s/L%d_fic%d.wav', basepath, i, j);
10                locutor_truth{index} = sprintf('L%d', i);
11                index = index+1;
12        end
13    end
14    for i=1:length(taux_reco)
15        if taux_reco(i) > 0.5
16            ↪ disp(sprintf('%s] Truth: %s - Predict:
17            ↪ %s (score: %d)', files{i}, locutor_truth{i},
18            ↪ 'locutor', taux_reco(i)));
19        else
20            ↪ disp(sprintf('%s] Truth: %s - Predict:
21            ↪ %s (score: %d)', files{i}, locutor_truth{i},
22            ↪ 'monde', taux_reco(i)));
23        end
24    end
25 end
```

En exécutant cette procédure pour le locuteur n°1, nous obtenons :

```
Locutor tested: L1
[REC0/L1_fic9.wav] Truth: L1 - Predict: locutor (score: 5.538741e-01)
[REC0/L1_fic10.wav] Truth: L1 - Predict: locutor (score: 5.898953e-01)
[REC0/L2_fic9.wav] Truth: L2 - Predict: monde (score: 3.611663e-01)
[REC0/L2_fic10.wav] Truth: L2 - Predict: monde (score: 3.409091e-01)
[REC0/L3_fic9.wav] Truth: L3 - Predict: monde (score: 2.352941e-01)
[REC0/L3_fic10.wav] Truth: L3 - Predict: monde (score: 2.533003e-01)
[REC0/L4_fic9.wav] Truth: L4 - Predict: monde (score: 2.595238e-01)
[REC0/L4_fic10.wav] Truth: L4 - Predict: monde (score: 4.242301e-01)
[REC0/L5_fic9.wav] Truth: L5 - Predict: monde (score: 3.606838e-01)
[REC0/L5_fic10.wav] Truth: L5 - Predict: monde (score: 3.946416e-01)
```

[REC0/L6_fic9.wav] Truth: L6 - Predict: monde (score: 3.036131e-01)
[REC0/L6_fic10.wav] Truth: L6 - Predict: monde (score: 3.117563e-01)
[REC0/L7_fic9.wav] Truth: L7 - Predict: monde (score: 4.783198e-01)
[REC0/L7_fic10.wav] Truth: L7 - Predict: monde (score: 4.789562e-01)
[REC0/L8_fic9.wav] Truth: L8 - Predict: monde (score: 2.603858e-01)
[REC0/L8_fic10.wav] Truth: L8 - Predict: monde (score: 3.055556e-01)
[REC0/L9_fic9.wav] Truth: L9 - Predict: monde (score: 3.669951e-01)
[REC0/L9_fic10.wav] Truth: L9 - Predict: monde (score: 3.439364e-01)
[REC0/L10_fic9.wav] Truth: L10 - Predict: monde (score: 1.780899e-01)
[REC0/L10_fic10.wav] Truth: L10 - Predict: monde (score: 2.306080e-01)

7 Fonction main

Voici la fonction `main` qui résume toutes les opérations effectuées.

```
1 % Nombre de fichiers d'apprentissage et de tests.
2 nbe_loc = 10;
3 nbe_fic = 8;
4 nbe_coef = 32; % number of cepstral coefficient,
5 threshold = 0.0001;
6 taille_fenetre = 512; % frame width,
7
8 %% Preprocessing
9 % Labelize data
10 etiquetage_total(taille_fenetre, nbe_loc, nbe_fic,
   ↪ threshold);
11
12 % Compute MFCC
13 parametrisation_total(taille_fenetre, nbe_loc,
   ↪ nbe_fic, nbe_coef);
14
15 %% Learning
16
17 % Constants
18 nb_gauss = 4; % number of gaussian,
19 locutor = 1; % locuteur id
20
21 % Affectation step
22 affectation('LABELS', 'MFCC', nbe_loc, nbe_fic,
   ↪ locutor);
23
24 % Learning step
25 [m, v, w, m_l, v_l, w_l] =
   ↪ apprentissage('MFCC/monde.mfcc',
   ↪ sprintf('MFCC/L%d.mfcc', locutor), nb_gauss);
26
27 % Test total
28 taux_reco = tests_total(taille_fenetre, nbe_coef, m,
   ↪ v, w, m_l, v_l, w_l);
29 disp(sprintf('Locutor tested: L%d', locutor))
30 print_truth_predict(taux_reco);
```

Pour aller encore plus loin je pense que nous aurions pû mettre en oeuvre une stratégie Grid Search pour trouver les meilleurs paramètres pour améliorer la précision de notre classification. Les 2 paramètres les plus importants sont

le nombre de coefficients cepstraux à retenir lors du calcul de la MFCC, et
le nombre de lois gaussiennes à utiliser lors de la phase d'apprentissage.