

Séparation de sources audio-visuelles : Formalisation et expérimentation

D. Sodoyer¹, L. Girin¹, C. Jutten², J.L. Schwartz¹

¹Institut de la Communication Parlée – INPG/Univ. Stendhal/CNRS

46, av Felix Viallet, 38031 Grenoble Cedex 1, France

Tél.: ++33 (0)476 57 47 12 - Fax: ++33 (0)476 57 47 10

Mél: sodoyer@icp.inpg.fr - http://www.icp.inpg.fr

²Laboratoire des Images et des Signaux – INPG/UJF/CNRS

ABSTRACT

In this paper, we present a new approach to the source separation problem in the case of multiple speech signals. The method is based on the use of automatic lipreading: the objective is to extract an acoustic speech input from other acoustic signals by exploiting its coherence with the speaker's lips movements. We consider the case of an additive stationary mixture. Firstly we present a theoretical framework showing that it is indeed possible to separate a source when some of its spectral characteristics are provided to the system. Then we address the case of audio-visual sources. We show how, if a statistical model of the joint probability of visual and spectral audio input is learnt to quantify the audio-visual coherence, separation can be achieved by maximising this probability.

1. INTRODUCTION

De nombreux travaux de recherche montrent l'existence d'une cohérence intrinsèque et même d'une complémentarité entre l'audition et la vision pour la perception de la parole [Sum87]. La contribution de la vision dans la perception de la parole bruitée permet d'obtenir un gain d'intelligibilité notable entre les conditions « audio seul » et « audio + visage du locuteur » [Ben94], [Rob98].

Dans leur étude, [Gir01] montrent qu'il est possible de réaliser un système de rehaussement de la parole dans un bruit blanc par filtrage utilisant l'image du locuteur. Dans la lignée de ce travail, nous cherchons maintenant à appliquer cette idée dans le cas d'un mélange de plusieurs sources de parole; le but étant de séparer l'une des sources. Ceci revient en fait à un problème de séparation de sources, problème qui a déjà suscité de nombreuses études en traitement du signal avec un certain nombre d'applications pour la parole. L'enjeu de ce travail est d'exploiter la cohérence entre l'image et le son de la parole. Nous aborderons ce principe en Section 2, les résultats étant exposés en Section 3. Puis en Section 4, nous présenterons quelques perspectives d'améliorations possibles afin de conclure.

2. MODELE

2.1 Présentation du problème

Nous considérons dans notre problème un système possédant N sources inconnues et localement stationnaires $s_i(t)$. La multiplication de ces dernières avec une matrice (P,N) de mélange A fournit P observations $x_k(t)$: $x = As$. L'opération de séparation consiste à estimer la matrice (N,P) de séparation B dont l'opération $y = Bx$ fournit N estimées $y_j(t)$ de $s_i(t)$ (figure 1). Nous prendrons le cas où $N = P$.

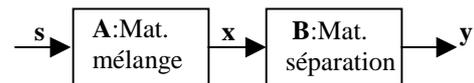


Figure 1: Schéma de la séparation de sources.

Dans les méthodes classiques d'ACI (Analyse en Composantes Indépendantes) la matrice de séparation B est estimée selon un critère de maximisation d'indépendance entre les sorties $y_j(t)$. Ici nous utilisons une information supplémentaire V_l correspondant aux dimensions géométriques des lèvres d'un locuteur fournies par la vidéo synchronisée avec son signal acoustique s_l , signal que nous chercherons à extraire d'un mélange de deux ou plusieurs sources (figure 2).

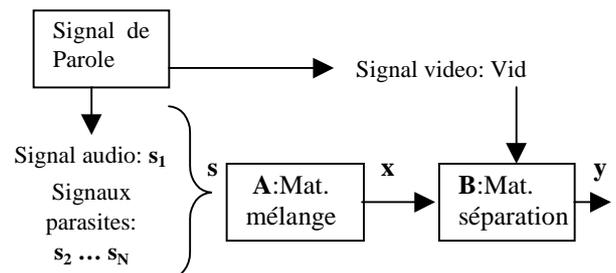


Figure 2 : Schéma de principe de la séparation de sources audio-visuelles.

2.2 Principe théorique

Dans notre étude, nous supposons simplement la décorrélation des sources, et une connaissance supplémentaire sur la source s_l que l'on désire extraire du mélange. Cette information supplémentaire correspond au signal visuel (associé à s_l) qui est approximativement lié partiellement au spectre de la fonction de transfert du conduit vocal, donc à

l'enveloppe spectrale du signal s_I . Nous pouvons donc associer à cette information spectrale, soit des coefficients d'autocorrélation, soit des coefficients d'énergie résultant d'un banc de filtre.

Critère d'autocorrélation

Nous supposons dans cette sous partie être capable de connaître une valeur d'autocorrélation du signal s_I pour un délai k donné. En normalisant cette valeur d'autocorrélation par l'énergie de s_I nous posons comme coefficient d'autocorrélation :

$$\gamma_k = \frac{R_{s_I s_I}(k)}{R_{s_I s_I}(0)} \quad (1)$$

Si, à la séparation, nous voulons que y_I soit une estimée de s_I , nous devons obtenir :

$$\frac{R_{y_I y_I}(k)}{R_{y_I y_I}(0)} = \gamma_k \quad (2)$$

Ainsi nous proposons de minimiser le critère :

$$f_{AC}(y) = C_k(y_I y_I)^2 \quad (3)$$

$$\text{avec : } C_k(y_i y_j) = R_{y_i y_j}(k) - \gamma_k R_{y_i y_j}(0) \quad (4)$$

Ce critère a la propriété d'être positif ou nul, et minimum (= 0) quand la séparation est achevée ($y_I = \lambda s_I$). Cependant nous devons vérifier que ce critère ne s'annule qu'à la séparation et non dans d'autres configurations. Pour ceci, nous introduisons la matrice $G=BA$ ce qui nous permet de poser :

$$y_p = \sum_n g_{pn} s_n \quad (5)$$

La question est de savoir si l'annulation du critère annule tous les g_{ln} pour $n \geq 2$. En introduisant (5) dans l'expression de l'autocorrélation nous obtenons :

$$R_{y_I y_I}(k) = \sum_{m,n} g_{Im} g_{In} R_{s_m s_n}(k) \quad (6)$$

Sous l'hypothèse que les sources sont décorréliées :

$$R_{y_I y_I}(k) = \sum_n g_{In}^2 R_{s_n s_n}(k) \quad (7)$$

En introduisant (7) dans (3), il faut pour annuler $f_{AC}(y)$:

$$C_k(y_I y_I) = \sum_n g_{In}^2 C_k(s_n s_n) = 0 \quad (8)$$

Par définition $C_k(s_I s_I) = 0$, ainsi :

$$C_k(y_I y_I) = \sum_{n \geq 2} g_{In}^2 C_k(s_n s_n) = 0 \quad (9)$$

Dans le cas de 2 sources, en supposant que $C_k(s_2 s_2)$ n'est pas nul, c'est à dire que les 2 sources n'ont pas la même propriété spectrale défini en (1), (9) implique que $g_{I2} = 0$ et ainsi nous obtenons $y_I = g_{I1} s_I$. Mais pour plus de 2 sources, l'éq. (9) n'assure pas l'annulation

des g_{ln} $n \geq 2$. Pour N sources il faut connaître $(N-1)$ valeurs de γ_k définies en (1) par $(N-1)$ délais k différents. Nous posons alors un nouveau critère $f_{AC}(y)$:

$$f_{AC}(y) = \sum_{k=1}^{N-1} C_k(y_I y_I)^2 \quad (10)$$

Ici $f_{AC}(y)$ est égal à zéro si et seulement si :

$$C_k(y_I y_I) = 0 \quad \forall k \in \{1 \dots (N-1)\} \quad (11)$$

Introduisant (9) dans (11) :

$$\begin{bmatrix} C_1(s_2 s_2) & \dots & C_1(s_n s_n) & \dots & C_1(s_N s_N) \\ \vdots & & \vdots & & \vdots \\ C_k(s_2 s_2) & & C_k(s_n s_n) & & C_k(s_N s_N) \\ \vdots & & \vdots & & \vdots \\ C_{N-1}(s_2 s_2) & & C_{N-1}(s_n s_n) & & C_{N-1}(s_N s_N) \end{bmatrix} \begin{bmatrix} g_{I2}^2 \\ \vdots \\ g_{In}^2 \\ \vdots \\ g_{IN}^2 \end{bmatrix} = 0 \quad (12)$$

Si la matrice $C\{C_k(s_n s_n)\}$ n'est pas singulière, cela implique que tous les g_{ln} de 2 à N sont nuls et ainsi $y_I = g_{I1} s_I$. La non-singularité traduit l'hypothèse que les sources ont des propriétés spectrales différentes.

Critère Spectral

De la même manière, nous pouvons supposer que nous connaissons un coefficient spectral fourni par un banc de K filtres passe-bandes. Posons $H_k(f)$ la réponse fréquentielle du filtre passe-bande RIF de canal spectral k ($1 \leq k \leq K$) et $h_k(t)$ sa réponse impulsionnelle. L'énergie de la source s_I à la sortie du processus de filtrage est fournie par l'autocorrélation de retard nul du signal filtré $h_k * s_I(t)$. Nous pouvons alors supposer connaître l'énergie normalisée du signal s_I dans un certain canal spectral:

$$\gamma_k = \frac{R_{(h_k * s_I)(h_k * s_I)}(0)}{R_{s_I s_I}(0)} \quad (13)$$

Comme dans le cas précédent, nous posons la fonction:

$$C'_k(y_i y_j) = R_{(h_k * y_i)(h_k * y_j)}(0) - \gamma_k R_{y_i y_j}(0) \quad (14)$$

Puis comme en (10), nous posons comme critère possible :

$$f_{SC}(y) = \sum_{k=1}^{N-1} C'_k(y_I y_I)^2 \quad (15)$$

Ce critère, fondé sur la connaissance de $(N-1)$ canaux spectraux, se comporte comme le précédent et conduit à la même équation (12).

En résumé, cette analyse théorique montre qu'il faut $(N-1)$ connaissances spectrales (coefficients d'autocorrélation ou coefficients spectraux) du signal s_I afin de pouvoir l'extraire d'un mélange additif à N sources. Notre étude s'est orientée sur l'extraction d'une seule source s_I . Pour la minimisation de nos critères [Sod01] nous avons utilisé la technique du gradient relatif introduisant la notion d'équivariance [Car96]. L'extraction d'autres sources supposera que l'on connaisse des informations spectrales du même

type pour chacune d'elles, ou dans le but d'une séparation totale, qu'il sera possible d'utiliser l'indépendance.

Dans notre application nous ne possédons en fait que des approximations de ces connaissances spectrales, celles-ci étant estimées par l'image des lèvres associées au signal de parole s_I . La section suivante a pour objet de faire le lien entre cette application et l'analyse théorique qui vient d'être vue.

2.3 Méthode Audio-visuelle

Sachant qu'il existe une relation statistique entre les dimensions labiales et le spectre du signal de parole d'un locuteur [Yeh98], nous construisons un modèle statistique audio-visuel $p(\mathbf{S}, \mathbf{V})$ comme étant une probabilité conjointe du vecteur visuel \mathbf{V} (contenant des dimensions labiales) et d'un vecteur \mathbf{S} représentant le spectre du signal audio. Nous définissons cette probabilité conjointe comme une somme de gaussiennes estimée par un algorithme Expectation-Maximisation (EM) itératif exploitant un corpus d'apprentissage afin de déterminer la moyenne et la matrice de covariance de chaque gaussienne.

L'algorithme de séparation consiste à trouver une matrice \mathbf{B} pour laquelle la sortie \mathbf{y}_I fournit un vecteur spectral le plus cohérent possible avec la vidéo \mathbf{V}_I associée à s_I . Il faut donc maximiser la probabilité audio-visuelle $p(\mathbf{Y}_I, \mathbf{V}_I)$ ou minimiser le critère:

$$f_{AV}(\mathbf{y}) = -\text{Log}(p(\mathbf{Y}_I, \mathbf{V}_I)) \quad (16)$$

Afin de faire le lien entre les eq. (15) et (17), considérons le cas où $p(\mathbf{Y}_I, \mathbf{V}_I)$ possède une seule gaussienne centrée

$$p(\mathbf{Y}_I, \mathbf{V}_I) = \frac{1}{\sqrt{(2\pi)^d \det(\mathbf{C})}} \exp\left[-\frac{1}{2} [\mathbf{Y}_I \mathbf{V}_I]^t \mathbf{C}^{-1} [\mathbf{Y}_I \mathbf{V}_I]\right] \quad (17)$$

où \mathbf{Y}_I est un vecteur spectral défini comme en Sect. 2.1 par K valeurs d'énergies fournies par un banc de filtre à K bandes passantes, \mathbf{V}_I est un vecteur contenant des dimensions labiales, et \mathbf{C} la matrice de covariance estimée avec le corpus d'apprentissage. Minimiser $f_{AV}(\mathbf{y})$, revient à minimiser le terme de l'exponentielle de (17). En décomposant \mathbf{C}^{-1} selon les termes impliquant les composantes auditives, visuelles et croisées, soit :

$$\mathbf{C}^{-1} = \begin{bmatrix} \mathbf{D} & \mathbf{E} \\ \mathbf{E}^t & \mathbf{F} \end{bmatrix} \quad (18)$$

nous obtenons un nouveau critère à minimiser :

$$f_{AV2}(\mathbf{y}) = \mathbf{Y}_I^t \mathbf{D} \mathbf{Y}_I + 2 \mathbf{V}_I^t \mathbf{E}^t \mathbf{Y}_I + \mathbf{V}_I^t \mathbf{F} \mathbf{V}_I \quad (19)$$

ou de manière équivalente:

$$f_{AV3}(\mathbf{y}) = (\mathbf{Y}_I - \mathbf{H} \mathbf{V}_I)^t \mathbf{D} (\mathbf{Y}_I - \mathbf{H} \mathbf{V}_I) \quad (20)$$

$$\text{avec } \mathbf{H} = -\mathbf{D}^{-1} \mathbf{E} \quad (21)$$

\mathbf{H} est alors la matrice de régression de \mathbf{V} vers \mathbf{S} . Ainsi nous faisons le lien avec le critère (15), la connaissance du vecteur \mathcal{Y}_k , contenant les coefficients γ'_k , étant remplacée par $\mathbf{H} \mathbf{V}_I$, et \mathbf{D} remplaçant une simple sommation dans (15) par une sommation plus complexe comprenant des rotations et des poids.

Etude du cas de 2 sources.

Sachant que nous cherchons à extraire la source s_I sur \mathbf{y}_I , dans le cas d'un mélange de 2 sources, et en imposant le terme $b_{11}=1$ de la matrice \mathbf{B} , nous nous intéressons alors à :

$$\mathbf{y}_I = \mathbf{x}_I + b_{12} \mathbf{x}_2 \quad \text{avec } b_{12} = \arg \min(f_{AV}(\mathbf{y})) \quad (22)$$

Quand \mathbf{y}_I est défini par la 1^{ère} eq. de (22), le spectre de \mathbf{Y}_I décrit une courbe dans l'espace spectral en fonction de b_{12} , la 2^{ème} eq. spécifiant la valeur optimale de b_{12} tel que \mathbf{y}_I soit le plus cohérent avec \mathbf{V}_I donc que \mathbf{y}_I soit proche de $g_I s_I$. Cependant, l'information vidéo \mathbf{V}_I peut être parfois associée à plusieurs spectres ce qui peut aboutir à une mauvaise séparation. C'est pourquoi, en nous permettant de considérer que les données audio et spectrales sont indépendantes trame à trame, nous spécifions la probabilité conjointe intégrée:

$$p(\mathbf{Y}_I(t, \dots, t-T), \mathbf{V}_I(t, \dots, t-T)) = \prod_{l=0}^T p(\mathbf{Y}_I(t-l), \mathbf{V}_I(t-l)) \quad (23)$$

En remplaçant dans l'eq. (16) la probabilité sur une trame par cette probabilité conjointe intégrée, nous pouvons espérer une meilleure estimation de la matrice \mathbf{B} .

3. RESULTATS EXPERIMENTAUX

3.1 Methodologie

Le corpus utilisé est un corpus mono locuteur composé des séquences $[\mathbf{V}_1 \mathbf{C} \mathbf{V}_2 \mathbf{C} \mathbf{V}_1]$ avec \mathbf{V}_1 et \mathbf{V}_2 une voyelle parmi [a, i, y, u] et \mathbf{C} une consonne parmi [p, t, k, b, d, g]. Il y a deux répétitions de chaque séquence, une pour l'apprentissage et l'autre pour le test.

Les signaux audio-visuels, synchrones, sont analysés par trame de 20 ms. Pour chaque trame, le système d'analyse labiale développé à l'ICP fournit deux paramètres labiaux (hauteur et largeur intero labiales) et une analyse spectrale fournit 32 canaux spectraux. Une analyse en composantes principales permet ensuite de réduire le nombre de composantes spectrales en perdant un minimum d'information. Par la suite, nous prendrons 5 ou 8 composantes principales, fournissant respectivement 85.5 % et 92.5 % de variance expliquée.

Le corpus d'apprentissage (96 stimuli de 24 trames en moyenne, soit à peu près 2300 trames au total) est traité par l'algorithme EM qui fournit une modélisation multi-gaussienne pour laquelle nous avons utilisé 6, 8 et 12 gaussiennes.

Critère d'évaluation

Nous considérons un mélange à deux sources :

$$\mathbf{x}_1 = a_{11} s_1 + a_{21} s_2 \quad ; \quad \mathbf{x}_2 = a_{21} s_1 + a_{22} s_2 \quad (24)$$

s_2 étant le signal perturbateur. Les sources étant normalisées en énergie, les RSBs définis pour s_1 sur chaque capteur x_i s'en déduisent comme :

$$RSB_{x_1} = 20\log(a_{11}^2/a_{12}^2) ; RSB_{x_2} = 20\log(a_{21}^2/a_{22}^2) \quad (25)$$

Et le RSB sur la sortie y_1 :

$$RSB_{y_1} = 20\log((a_{11}+b_{12}a_{12})^2/(a_{12}+b_{12}a_{22})^2) \quad (26)$$

Nous avons pris pour la suite le mélange suivant :

$$a_{11}=2 \quad a_{12}=1 \quad a_{21}=3 \quad a_{22}=5 \quad (27)$$

Ceci fournit théoriquement une solution $b_{12}=-0.2$, et des valeurs de RSB_{x_i} respectivement de 6 et -4.4 dB. L'expérimentation a été faite avec 96 stimuli de test (cf Sect. 3.1) contenant à peu près 2300 trames. La procédure de séparation est la suivante : y_1 étant défini par la 1^{ère} partie de (22), on estime son spectre Y_1 selon le processus décrit en Sect. 3.1, puis on estime b_{12} selon la 2^{ème} partie de (22). Une première optimisation a été effectuée avec l'algorithme du simplexe [Kli00].

3.2 Résultats

Sur les 2300 trames de test, nous avons comptabilisé le nombre de trames pour lesquelles la valeur de b_{12} estimée appartenait à l'intervalle $[-0.15, -0.25]$ centrée sur la solution, trames pour lesquelles RSB_{y_1} est supérieur à 14 dB donc correctement débruitées. Nous présentons sur la figure 3 le pourcentage de trames correctement débruitées, pour une modélisation à 12 gaussiennes, 5 ou 8 composantes principales, et un nombre de trames T pour l'intégration temporelle (23) variant de 1 à 20.

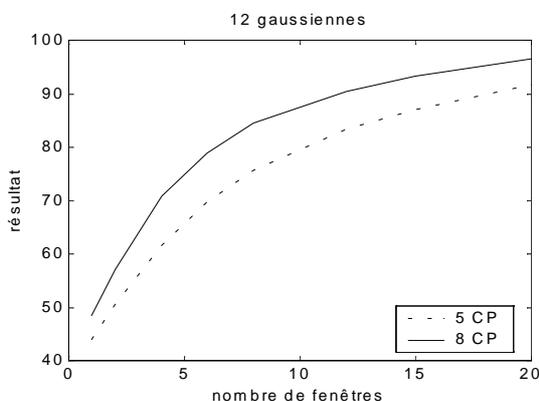


Figure 3 : Pourcentage de l'estimation correcte de b_{12} sur le corpus de test en fonction de T .

Au vu de cette figure il est assez facile de se rendre compte de l'importance de l'intégration temporelle. Pour $T=20$, 96 % de trames sont correctement débruitées si l'on utilise 8 composantes principales. L'augmentation du nombre de composantes principales et de gaussiennes pour la modélisation accroît légèrement les performances.

4. CONCLUSION

Il apparaît que cet algorithme audio-visuel possède les caractéristiques d'être théoriquement cohérent et techniquement réalisable. Cependant il aura besoin de

plusieurs améliorations fondamentales. D'abord nous tenterons d'appliquer la méthode du gradient relatif vue en Sect. 2.1 à la place de celle du simplexe afin d'améliorer l'optimisation. Deuxièmement notre modèle statistique de probabilité conjointe audio-visuelle reposant sur une hypothèse d'indépendance trame à trame pourrait être remplacé par une probabilité plus complexe introduisant par exemple des modèles de Markov. D'autre part, si l'on compare notre algorithme à d'autres algorithmes classiques de séparation de sources, nous avons des performances semblables en termes de RSB sur la première source extraite [Kli00]. Il est vrai que de tels algorithmes permettent de séparer toutes les sources mais sans pouvoir contrôler et identifier une permutation possible entre les sources et les signaux de sortie de trame à trame. Ainsi notre algorithme dispose ici d'un atout non négligeable dans le cas de signaux audio en étant capable de combler cette carence. Enfin, un objectif à terme est bien sûr d'utiliser conjointement critère audio-visuel et critère d'indépendance. Ceci devrait s'avérer très utile dans des cas plus complexes impliquant des mélanges convolutifs ou des mélanges possédant moins de capteurs que de sources.

BIBLIOGRAPHIE

- [Sum87] Summerfield Q. (1987), in *Hearing by Eye : The Psychology of Lipreading*, edited by Dodd B. and Campbell R (Lawrence Erlbaum Associates, London), pp. 3-51.
- [Ben94] Benoît C. Mohamadi T. Kandel S. (1994) "Effects of phonetic context on audio-visual intelligibility of French", *JSHR*, Vol. 37, pp. 1195-1293.
- [Rob98] Robert-Ribes J.Schwartz J.L. Lallouache T. Escudier P. (1998), "Complementary and synergy in bimodal speech : auditory, visual, and audio-visual identification of French oral vowels in noise", *JASA*, Vol. 103 :6, pp. 3677-3689.
- [Gir01] Girin L. Schwartz J.L. & Feng G. (2001) "Audio-Visual enhancement of speech in noise", *JASA*, Vol. 109:6, pp. 3007-3020.
- [Com95] Comon P. (1995), "Independent Component Analysis, a new concept ?", *Sig. Proc.*, Vol. 36:3, pp. 287-314.
- [Sod01] Sodoyer D. (2001), "Séparation de sources Audio-visuelles: de la formalisation à l'expérimentation", Master, INP Grenoble.
- [Car96] Cardoso J.F. & Laheld B. (1996), "Equivariant adaptive source separation", *IEEE Trans. on SP*, Vol. 44:12, pp. 3017-3030.
- [Yeh98] Yehia H. Rubin P. & Vatikiotis-Bateson (1998), "Quantitative association of vocal-tract and facial behavior", *Speech Comm.*, Vol. 26, pp. 23-43.
- [Kli00] Klinkisch J. (2000), "Séparation audio-visuelle pour la reconnaissance de chiffre dans le bruit", Master, INP Grenoble.