

Comparaison de trois outils de détection automatique de prééminences en français parlé

¹N. Obin, ^{2,5}J.-Ph. Goldman, ³M. Avanzi, ⁴A. Lacheret

¹IRCAM, ²Université de Neuchâtel et Université de Paris X Nanterre, ³Université de Genève,

⁴Université de Paris X, MODYCO Nanterre et IUF, ⁵Université de Louvain-la-Neuve

Nicolas.Obin@ircam.fr; Jean-Philippe.Goldman@lettres.unige.ch; Mathieu.Avanzi@unine.ch; anne@lacheret.com

ABSTRACT

This paper presents the inner details of three different algorithms for prominence detection. On the basis of a 50-minute corpus made of five speaking styles and manually annotated for prominence, a quantitative evaluation compares the three approaches.

Keywords: spontaneous speech, prosody, prominence, automatic detection.

1. INTRODUCTION

L'annotation des prééminences accentuelles qui scandent le discours spontané ne peut plus être envisagée manuellement, étant donné d'une part l'aspect subjectif de la démarche, d'autre part le temps de codage demandé, d'autant plus coûteux qu'on souhaite aboutir à une annotation rigoureuse et stabilisée à l'issue de la confrontation de plusieurs expertises de codage. Une telle pratique manuelle est d'autant moins concevable qu'il s'agit de brasser des données de plus en plus volumineuses, quantitativement représentatives pour pouvoir les exploiter de manière fiable sous l'angle de l'analyse structurale et fonctionnelle. En pratique, les méthodes de traitement automatiques se développent [1] [2] et doivent continuer à se développer pour prendre le relais du codage manuel et, tout en capitalisant les connaissances acquises, les faire évoluer et enrichir les analyses prosodiques de la parole non lue, des discours formels aux discussions à bâtons rompus. Cette communication s'inscrit dans cette problématique. Un corpus de français parlé échantillonné en différents genres, synthétiquement présenté dans la section 2, est utilisé pour comparer trois outils de détection automatique de prééminences, exposés dans la section 3. La section 4 fait état du protocole d'évaluation et la partie 5 présente les résultats.

2. MATERIEL D'ETUDE

2.1. Le corpus

Pour cette étude, nous avons utilisé un corpus échantillonné en cinq différents genres et styles de parole : discours politique (DP), descriptions d'itinéraires (ITI), récits de vie (RCV), journaux radiophoniques (JP) et interviews radiophoniques (IRT). Les locuteurs sont tous des francophones natifs, et proviennent de France, de Suisse ou de Belgique). On en trouvera une présentation exhaustive dans [3].

2.2. Prétraitement et transcription

Le corpus a été transcrit et aligné semi-automatiquement (phonèmes, syllabes et mots graphémiques) sous Praat [4] avec EasyAlign [5]. Les alignements ont fait l'objet d'une vérification manuelle par deux experts humains, qui ont ensuite consigné dans une tire dédiée (tire *delivery*) les informations relatives au statut des syllabes (syllabes prééminentes ou non ; syllabes particulières contenant des phénomènes propres au langage spontané, *i.e.* hésitations, faux-départ, schwas post-toniques, bruits de bouche, etc). On se reportera aux travaux de [6] pour un compte-rendu exhaustif sur la procédure et ses origines. Au total, le corpus est composé de 12851 intervalles syllabiques. 973 d'entre eux ont été exclus via la tire *delivery*, 3244 syllabes ont été annotées prééminentes, 8634 ont été codées ni prééminentes ni associées à un marqueur *delivery*, soit 11878 à traiter automatiquement.

3. LES LOGICIELS DE DETECTION

Dans cette partie, nous décrivons les principes de base de chacun des trois outils dont nous souhaitons comparer les performances sur une tâche de détection automatique de prééminences.

3.1. Anolor

Anolor est un logiciel d'analyse implémenté sous Matlab [7] qui fonctionne avec des fichiers xml en sortie de Praat. Initialement, le logiciel a été conçu dans le but de faire émerger des critères acoustiques robustes en vue de segmenter automatiquement le continuum sonore en unités d'intégration prosodique maximales, ou **périodes intonatives**. Un ensemble de fonctions récemment élaboré permet aujourd'hui une détection automatique des syllabes prééminentes à l'intérieur des périodes identifiées par le logiciel. Cette procédure a été présentée et expérimentée pour la première fois dans [8]. Elle repose sur la mise au jour des variations significatives de hauteur et de durée par rapport à la moyenne globale de l'ensemble des syllabes qui composent une période. Appelons $M_h(P)$ la moyenne et $E_h(P)$ l'écart-type de la fondamentale F0 sur une période P. Une syllabe s de P est dite prééminente pour la hauteur si elle contient un maximum local de la F0, noté h(s), vérifiant la condition :

$$h(s) > M_h(P) + K_h * E_h(P)$$

où K_h est un paramètre ajustable indépendant du locuteur fixé à 1,5 par défaut.

Le calcul est le même pour la durée (pour une justification de ce seuil *a priori*, voir [9] et [10]). Les pauses silencieuses ont été également utilisées pour affiner la

détection des syllabes proéminentes : toute syllabe suivie d'une pause est identifiée comme proéminente, *i.e.* saillante perceptivement.

3.2. ProsoProm

La détection automatique des proéminences par l'outil ProsoProm se déroule en trois étapes successives : 1. segmentation et stylisation des noyaux vocaliques, 2. extraction et relativisation des paramètres acoustiques et 3. décision du statut proéminent pour chaque syllabe.

La première étape utilise une version adaptée du ProsoGram de Mertens [11], développé pour de la transcription prosodique semi-automatique, pour repérer et styliser le noyau vocalique de chaque syllabe (le noyau étant délimité comme la portion voisée qui « a suffisamment d'intensité », en se basant sur des seuils d'intensité relativement au maximum d'intensité local). Cette étape permet d'éliminer des erreurs de détection de F0 aux frontières de voisement. Puis, la courbe mélodique de ce noyau est stylisée en un ou plusieurs segments, plats ou avec une pente mélodique, selon des paramètres perceptuels comme le glissando. L'adaptation de ProsoGram [12] permet d'augmenter le nombre de noyaux effectivement stylisés et d'affiner certains paramètres de réglages pour cette application spécifique. Dans un second temps, plusieurs paramètres acoustiques sont estimés pour chaque syllabe, à savoir :

1. la durée syllabique, préférée à la durée du noyau car ce dernier est fortement contraint par le voisement de chacun des phonèmes composant la syllabe
2. la hauteur mélodique maximale du noyau stylisé, en demi-tons, considéré comme cible mélodique atteinte
3. l'amplitude du mouvement montant, en demi-tons. Le mouvement descendant est écarté selon l'idée qu'il est plutôt, à l'exception peut-être des syllabes finales, la manifestation d'un relâchement articulo-articulaire.
4. la durée de la pause subséquente en secondes. Ceci restreint fortement l'étude à des langues oxytoniques comme le français.

Les deux premiers paramètres sont « relativisés » par rapport aux noyaux adjacents (des deux syllabes précédentes et de la syllabe suivante). Ceci permet de rendre ces paramètres indépendants du débit et du registre de parole. L'empan de relativisation choisi est volontairement très local rejoignant l'hypothèse que malgré l'existence évidente d'unités intonatives plus larges, l'appui sur les syllabes immédiatement adjacentes est primordial pour réaliser un effet de contraste.

Finalement, la stratégie de décision consiste à comparer chacun des paramètres avec un seuil et de déclarer proéminente toute syllabe dont un des paramètres est au-dessus de son seuil propre. ProsoGram permet également de produire une sortie graphique dans laquelle les segmentations phonétiques et lexicales sont associées aux au tracé des noyaux stylisés, dont la couleur dépend de leur caractère proéminent. Son fort pouvoir explicatif permet un diagnostic qualitatif des erreurs subsistantes. Pour cette étude, une optimisation des quatre seuils a été faite par une approche dichotomique pour chaque sous-corpus d'entraînement.

3.3. IrcamProm

IrcamProm est un outil de détection de proéminences développé dans le cadre d'*IrcamCorpusTool* [13]. Il repose sur un modèle statistique de la proéminence suivant un protocole décrit en [14]. Ce modèle se décompose schématiquement en trois étapes :

1. Une étape d'extraction des paramètres acoustiques de la parole : paramètres acoustiques (hauteur, durée, intensité, information spectrale), mesures réalisées sur ces paramètres (statique, dynamique, informations de forme) et horizon temporel de relativisation (aucun, syllabes adjacentes, groupe accentuel, groupe prosodique).
2. Une étape de sélection des corrélats acoustiques de la proéminence, à partir d'un algorithme de sélection de descripteurs qui détermine l'ensemble des paramètres acoustiques qui permettent la meilleure discrimination des syllabes proéminentes et non proéminentes.
3. Une étape de modélisation, réalisée par un modèle de mélange de gaussiennes pour les syllabes proéminentes et non-proéminentes. Une fois les paramètres du modèle estimés, la décision de proéminence est réalisée suivant le critère de maximum a posteriori.

Les paramètres acoustiques retenus pour cette étude ont été obtenus par optimisation sur le corpus de parole lue présenté en [14]. Ils consistent en un jeu de paramètres acoustiques faisant intervenir plusieurs types d'information (hauteur, durée, information spectrale) relativisés sur plusieurs horizons temporels (absence, syllabes adjacentes et groupe prosodique). Nous les présentons dans l'ordre de leur importance relative pour la modélisation de la proéminence: la durée de la syllabe ; la valeur moyenne de la hauteur sur la syllabe relativisée par rapport à sa valeur moyenne sur les syllabes adjacentes ; la durée du noyau; la valeur moyenne de la sonie dans la 1^{ère} bande de Bark sur la syllabe relativisée par rapport à sa valeur sur la syllabe suivante ; la valeur moyenne de la sonie dans la 18^{ème} bande de Bark relativisée par rapport à sa valeur sur la syllabe précédente ; la valeur minimum du débit local sur la syllabe relativisée par rapport à celle du débit local sur le groupe prosodique ; la courbure du débit local sur la syllabe. Ce modèle a été élaboré sur un corpus de parole monolocuteur lue. Nous supposons dans cette étude que ces paramètres estimés sur un corpus particulier demeurent inchangés pour le corpus présentement étudié, hypothèse qu'il restera à vérifier. En revanche, pour adapter notre modèle à une détection de proéminences sur de la parole spontanée multi-locuteur dans divers genres discursifs, nous avons normalisé l'ensemble des paramètres par rapport à leur valeur moyenne et déviation standard sur chaque groupe prosodique. Cette « relativisation » est pratiquée afin de normaliser l'influence du locuteur et du genre de discours sur la distribution des paramètres acoustiques.

4. PROTOCOLE D'EVALUATION

Les outils présentés ont été comparés sur le corpus décrit suivant un protocole de *validation croisée*. Ainsi, le corpus a été préalablement échantillonné suivant les 5 genres de discours. Les outils sont alors entraînés sur 4

d'entre eux, puis validés sur le genre non observé. La procédure est réitérée en permutant les genres de discours utilisés pour l'apprentissage et le corpus utilisé pour la validation. Une telle méthode présente trois avantages : i) elle permet de tester la *capacité de généralisation* des outils en les testant sur des données non-observées au cours de l'apprentissage, ii) le choix de la validation sur un genre de discours non-observé et *a priori* suffisamment indépendant des genres observés permet de tester la *robustesse* des méthodes, et iii) les changements de corpus d'entraînement et de validation permettent de tester la *sensibilité* de la performance par rapport au corpus d'entraînement. L'outil Analor étant un système de décision à base de règles, et par conséquent non susceptible d'être entraîné, nous avons utilisé les paramètres standard de cet outil directement sur l'ensemble des genres discursifs. Ce biais inséré est contrebalancé par le fait que les paramètres de cet outil n'ont pas été réglés sur le corpus que nous traitons présentement. La mesure de performance des outils a été choisie comme étant la *F-mesure* de la classe *proéminence*, préférée à la mesure traditionnelle de *précision globale* ([1], [2]). Ce choix se justifie de la manière suivante : nous considérons d'abord que seules les erreurs sur l'estimation de la proéminence nous intéressent. Or, la mesure de *précision globale* prend en compte l'ensemble des classes (moyenne des erreurs sur chaque classe pondérées par le nombre d'observations de ces classes respectives). Et dans le cas du phénomène de proéminence, nous savons que les syllabes non-proéminentes sont beaucoup plus nombreuses que les syllabes proéminentes, ce a pour effet de favoriser la bonne estimation des syllabes non-proéminentes et de cacher partiellement les résultats obtenus pour les syllabes proéminentes. Cette mesure montre généralement des scores flatteurs, et qui ne sont pas nécessairement instructifs sur la qualité de la détection de la proéminence. La *f-mesure* de classe proéminence uniquement permet de résoudre ce problème et, de plus, de faire la synthèse des performances en terme d'insertion et de délétion des proéminences. La comparaison de la performance des outils est suivie d'une analyse automatique des différences de comportement observées entre eux. Cette analyse est fondée sur le *test de MacNemar* [15] qui permet de comparer les différences observées sur les erreurs de classification entre les outils et de diagnostiquer si ces différences sont statistiquement significatives.

5. RESULTATS ET DISCUSSION

Nous présentons dans la Table 1 les performances des outils sur chaque genre de discours. Nous voyons de manière générale que les performances se situent aux environs de 70% de *F-mesure* en s'échelonnant de 70% à 75%. Pour permettre une comparaison avec les méthodes de référence, les précisions globales associées sont respectivement de 84.9% pour Analor, 84.2% pour ProsoProm et de 86.2% pour IrcamProm. Ce niveau de performance est comparable aux performances observées sur des méthodes de référence (80-88%), dans une tâche

cependant plus difficile, tant au niveau du corpus étudié (parole spontanée) que du protocole d'évaluation (sur un genre de discours non observé et dont les structures de proéminence sont susceptibles de différer des genres de discours observés).

	iti	irt	rcv	dp	jp	Total
Analor	68,7	71,5	63,2	74,5	70,8	69,7
ProsoProm	74,4	73,4	71,0	74,9	72,5	73,2
IrcamProm	75,8	74,3	76,3	76,00	75,0	75,4

Table 1 : *F-mesure* des outils de détection sur chaque genre de discours **du corpus d'étude**.

IrcamProm présente la meilleure performance globale et par genre de discours. Analor apparaît plus sensible au genre de discours que les deux autres outils, avec une déviation standard relative 6 fois plus élevée que les outils ProsoProm et IrcamProm. Ceci montre à la fois les limites de capacité de généralisation d'un système à base de règles et sans adaptation au locuteur ni au genre de discours (le seuil proposé étant fixé).

Une analyse approfondie montre des tendances analogues pour l'ensemble des outils. Les trois logiciels présentent un minimum de performance sur un fichier de demande d'itinéraire. Cette propriété suggère : 1. la structure d'une proéminence diffère au moins quantitativement selon les genres de discours, et/ou 2. un contraste proéminence/non proéminence plus ou moins marqué selon les genres.

Un diagnostic des erreurs montre des comportements différents suivant les outils : Analor présente une tendance à la délétion de proéminence (rappel = 63.6%, précision = 77.2%) alors que ProsoProm présente une tendance à l'insertion de proéminence (rappel = 78.9%, préc. = 68.2%). IrcamProm présente un compromis entre délétion et insertion (rappel = 76.4%, préc. = 74.5%).

La différence de performance entre l'outil Analor et les outils ProsoProm et IrcamProm est doublement instructive : ainsi, si les performances de ProsoProm et de IrcamProm se révèlent robustes indépendamment du genre de discours, on observe une baisse sensible des performances d'Analor sur les genres de descriptions d'itinéraires et de récit de vie (parole spontanée), ce qui suggère un contraste de proéminences qui diffère des autres genres de discours. Si nous associons cette observation à la tendance de délétion mentionnée, nous pouvons dès lors supposer que la parole spontanée se manifeste par un moindre contraste de proéminence. Enfin, une analyse du test de MacNemar ($p = 0.005$) sur l'ensemble du corpus met à jour une différence significative entre l'outil IrcamProm et les deux autres outils, mais pas entre ces deux derniers. Néanmoins, cette estimation globale doit être relativisée par le fait que cette différence est variable suivant les genres (nous observons une absence de diagnostic de différence sur les corpus d'itinéraires et de journaux parlés).

6. CONCLUSION

Nous avons exposé les performances de trois outils de détection automatique de proéminences dans un contexte de traitement de corpus de français parlé. Ces outils présentent chacun des particularités qui expliquent sans

doute les scores obtenus et qui nous amènent à nous interroger et à mieux préciser l'impact des principes sous-jacents à la détection dans chaque système. Le premier, Analor, est un système à base de règles, alors que les deux autres sont fondés sur l'apprentissage. ProsoProm repose sur une stylisation perceptive de la courbe mélodique alors qu'Analor et IrcamProm se fondent sur des paramètres acoustiques bruts. IrcamProm prend en compte un nombre important de paramètres, alors que les deux autres algorithmes se focalisent sur la durée et les variations de F0, au mieux l'intensité (ProsoProm). Dernière différence : Analor travaille sur une fenêtre périodique tandis que ProsoProm et IrcamProm n'intègrent que les contextes syllabiques immédiats.

A partir de ces constats, plusieurs questions se posent. et sont à prendre en compte dans la suite de ce travail, à savoir le diagnostic quantifié des erreurs, en particulier :

- Peut-on mesurer l'impact de la stylisation perceptive (ProsoProm) par rapport à des approches travaillant uniquement sur la courbe acoustique brute ?
 - Comment expliquer le paradoxe apparent : pourquoi le système à bases de règles (Analor) est plus sensible aux variations de genre discursif que les autres ?
 - Quel est le meilleur compromis à trouver entre les résultats obtenus (performance de l'outil) et la complexité algorithmique demandée (voir IrcamProm) ?
 - Quel est la meilleure fenêtre à prendre en compte dans un système par apprentissage ? Dans quelle mesure, dans un tel système, la période comme fenêtre d'analyse améliore-t-elle ou non les résultats par rapport à la seule prise en compte des syllabes immédiatement adjacentes ?
- Enfin, rappelons que l'évaluation des outils s'est déroulée sur des syllabes non-marquées *delivery*, laissant ainsi de côté ces phénomènes propres à l'oral, et pour lesquels il est désormais possible d'envisager une détection automatique. Voilà l'ensemble des questions qui articulent notre programme de recherche à venir.

7. REMERCIEMENTS

Trois cadres institutionnels portent ce projet :

- Fonds National de la recherche scientifique Suisse (subside n°100012-113726/1, "La structure interne des périodes", Université de Neuchâtel),
- ANR Rhapsodie 07 Corp-030-01, Corpus prosodique de référence du français parlé.
- Programme Wist2 Convention n°616422, financé par la Région wallonne (Belgique) Projet *EXPRESSIVE*

Nous souhaitons remercier ici A.-C. Simon et F. Poiré, initiateurs de la méthode de codage manuel utilisée ; ainsi que B. Victorri pour sa collaboration active dans les différentes étapes de la modélisation.

BIBLIOGRAPHIE

- [1] F. Tamburini, & C. Caini, An Automatic System for Detecting Prosodic Prominence in American English

Continuous, *Speech International Journal of Speech Technology* 8: 33-44, 2005.

- [2] A. Rosenberg and J. Hirschberg, "Detecting pitch accent using pitch corrected energy-based predictors," *Interspeech'07*, 2777-2780, 2007.
- [3] A.-C. Simon, M. Avanzi & J.-P. Goldman. La détection des prééminences syllabiques. Un aller-retour entre l'annotation manuelle et le traitement automatique. Article soumis au *1^{er} Congrès Mondial de Linguistique Française*, 2008.
- [4] P. Boersma & D. Weenink. Praat: doing phonetics by computer (Version 4.5). www.praat.org, 2008.
- [5] J.-Ph. Goldman (2008). EasyAlign: a semi-automatic phonetic alignment tool under Praat, Avril 2008, <http://latlcui.unige.ch/phonetique/easyalign>.
- [6] M. Avanzi, J.P. Goldman, A. Lacheret-Dujour, A.C. Simon & A. Auchlin. Méthodologie et algorithmes pour la détection automatique des syllabes prééminentes dans les corpus de français parlé. *Cahiers of French Language Studies*, 13/2, 2007.
- [7] A. Lacheret-Dujour & B. Victorri. La période intonative comme unité d'analyse pour l'étude du français parlé : modélisation prosodique et enjeux linguistiques. *Verbum*, 24/1-2. 55-73, 2002.
- [8] M. Avanzi, A. Lacheret-Dujour & B. Victorri. ANALOR. A Tool for Semi-Automatic Annotation of French Prosodic Structure, *Speech Prosody*, 2008.
- [9] M. Rossi & al. : L'intonation, de l'acoustique à la sémantique, Paris, Klincksieck, 1981.
- [10] A. Lacheret-Dujour & F. Beaugendre : La prosodie du français, Paris, CNRS éd., 1999.
- [11] Mertens, P. The Prosogram: Semi-Automatic Transcription of Prosody based on a Tonal Perception Model", in B. Bel & I. Marlien (eds.) *Speech Prosody*, Nara (Japan), 2004
- [12] J.-P. Goldman, M. Avanzi, A. Lacheret-Dujour, A.-C. Simon & A. Auchlin. A Methodology for the Automatic Detection of Perceived Prominent Syllables in Spoken French., *Interspeech'07*, pp. 91-120, 2007.
- [13] C. Veaux, G. Beller, D. Schwarz, & X. Rodet, Ircamcorpustools: an extensible platform for speech corpora exploitation, *ELREC'08*, Marrakech, 2008.
- [14] N. Obin, X. Rodet & A. Lacheret-Dujour. French Prominence: a probabilistic framework. *ICASSP'08*, Las Vegas, Nevada, USA, 2008.
- [15] T.G. Dietterich, Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms, in *Neural Computation*, Vol. 10, p. 1895-1923, 1998