

Extracting the acoustic features of interruption points using non-lexical prosodic analysis

Matthew P. Aylett

ICSI, UC Berkeley, USA and
CSTR, University of Edinburgh, UK

Abstract

Non-lexical prosodic analysis is our term for the process of extracting prosodic structure from a speech waveform without reference to the lexical contents of the speech. It has been shown that human subjects are able to perceive prosodic structure within speech without lexical cues. There is some evidence that this extends to the perception of disfluency, for example, the detection of interruption points (IPs) in low pass filtered speech samples. In this paper, we apply non-lexical prosodic analysis to a corpus of data collected for a speaker in a multi-person meeting environment. We show how non-lexical prosodic analysis can help structure corpus data of this kind, and reinforce previous findings that non-lexical acoustic cues can help detect IPs. These cues can be described by changes in amplitude and f_0 after the IP and they can be related to the acoustic characteristics of hyper-articulated speech.

1. Introduction

Human subjects respond to prosodic structure without necessarily understanding the lexical items which make up the utterance. For example event-related brain potential (ERP) studies have shown a reliable correlation with phrase boundaries when utterances are made lexical nonsensical, either by humming the words, or by replacing them with nonsense words [9]. The use of prosodically rich pseudo speech for artistic purposes (such as R2D2 in star wars, and The Teletubbies amongst others) reinforce these findings. This effect, of apparently understanding prosodic structure without lexical cues, extends to the human perception of disfluency. Lickley [7] showed that human subjects could recognise interruption points, the boundary between disfluent and fluent speech, in low pass filtered speech where no lexical cues were present.

Non-lexical prosodic analysis (NLPA) attempts to mimic this human ability of non-lexical prosodic recognition. Initially, interest in NLPA was motivated largely by the objective of improving automatic speech recognition (ASR) technology, for example, by pre-processing the speech to find syllables [5] or prosodic prominence [3]. However, improvements in statistical modeling in ASR meant that, often, the speech recogniser itself was best left to model prosodic effects internally. Recently, there has been a renewed interest in NLPA techniques in order to address the problem of recognising, segmenting, and characterising very large spontaneous speech databases. Tamburini and Caini [10] point out that identifying prosodic phenomena is useful, not only for ASR and speech synthesis modeling, but also for disambiguating natural language and for the construction of large annotated resources. In these cases, the ability to recognise prosodic structure without lexical cues has two main advantages:

1. It does not require the resource intensive, and language dependent, engineering required for full speech recognition systems.
2. It can offer a means of modeling the human recognition of prosodic structure which in turn could lead to an improved understanding of human speech perception and production.

The ability of human subjects to recognise interruption points (IPs) without lexical information raises the question of whether NLPA can do as good a job. Although previous work has looked at this problem in some depth (e.g. [4], [7]), NLPA offers the prospect of a structured analysis that could be carried out automatically over very large speech databases. In addition, the presence of previous detailed studies allows us to validate the overall approach.

The non-lexical detection of IPs is also of interest from the perspective of determining dialogue structure. Recent work suggests that disfluency patterns could be used to signal the speakers' cognitive load [1] and thus might be used to determine areas in dialogue involving complex concepts, ideas or planning.

We will first describe in more detail the corpus of speech we analysed and the IP phenomena. Next, we will present the details of the NLPA we applied to this corpus followed by results for a set of acoustic features which may cue the non-lexical perception of IPs. Finally, we will discuss limitations with the approach and possible future work.

2. Corpus and disfluency coding

Our data was selected from the ICSI meeting corpus [6]. This consists of 75 dialogues collected from the regular weekly meetings of various ICSI research teams. Meetings in general run for under an hour and have on average 6.5 participants each recorded on a separate acoustic channel. The speech is segmented into spurts, defined as periods of speech which have no pauses greater than 0.5 seconds.

The data we present here is taken from a single speaker¹ taken from two dialogues. Disfluencies are coded as part of the dialogue act coding [2], where interruption points are shown as a hyphen in the speech transcription. In order to avoid complexity caused by multi-speaker interaction and multiple disfluencies, we looked only at phrase boundaries and IPs where:

- The same speaker continued speaking after the interruption point or phrase break
- No other speakers were speaking within 0.5 seconds of the break
- There was at least 0.5 seconds between any breaks.

Pause duration is the clearest acoustic cue of a prosodic break and can be used to disambiguate between IPs and phrase boundaries with some success. In general, the longer the pause, the more likely the break is a phrase boundary. However there are plenty of examples of phrase boundaries followed by a short pause. An interesting question is whether

¹We hope to increase the scale of this analysis for the final version of the paper

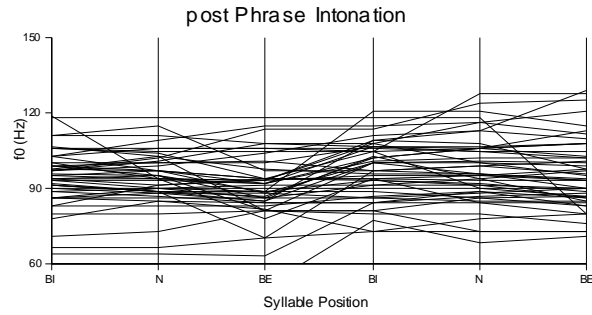
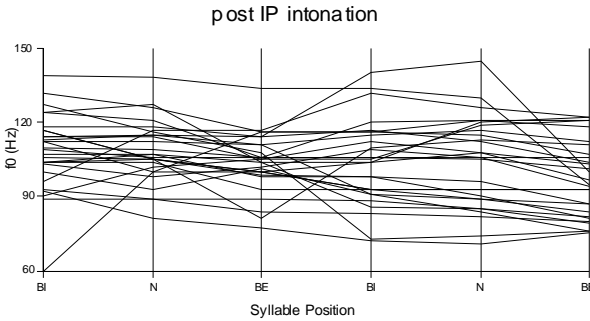


Figure 2: F0 across the two syllable nuclei following both IPs and phrase breaks where no or minimal pause cues are present. (BI - initial boundary of syllable nucleus, N - centre of syllable nucleus, BE - end boundary of syllable nucleus).

applied to reduce recognition error from 25% down to 16%. However, the extent to which these utterances contained IPs was not reported. In Lickley [7], human judgments of low pass filtered speech utterances show a significant, although far from consistent, effect across materials. Human subjects tended to misclassify disfluent utterances as fluent utterances more than visa-versa with the best group of human subjects correctly classifying 34% of disfluent utterances as disfluent. In addition, the significant effect in this study appeared to be dominated by the presence and differences of pause durations rather than other acoustic cues.

In this data, as stated earlier, only boundaries with pauses not discernable to the autosegmenter where examined. We compared the results for IPs and normal phrase breaks. As in [7] [4], we looked at acoustic cues in the form of f0 variation, syllabic nucleus amplitude and syllabic nucleus duration after the boundary point.

5. Results

We began by looking at the f0 change across the two syllables to the right of the boundaries. Shown in figure 2 are six f0 points. These values are taken from the first two syllable nuclei found with NLPA subsequent to the phrase or IP boundary. It is interesting to note a lack of a homogeneous f0 structure in either IP or for PH (Phrase conditions). However, differences are clearly present between both groups. F0 in the IP case tends to be higher and varies more throughout the two syllables.

On the basis of this plot we chose three f0 features to examine statistically: the f0 before the boundary, the f0 following the boundary and the variance of the f0 across the two syllables following the boundary. In addition, we combined the log of the raw amplitude of the first following syllable with the log of the duration of its nucleus by multiplying the factors together to give an overall prominence factor. Thus short, high energy syllable nuclei where regarded as having similar prominence to long, lower energy syllable

Table 1: Independent t-test for acoustic cues following IPs and Phrase Boundaries.

Acoustic Feature	t	df	Sig. (2-tailed)	Bonferroni correction
f0 pre boundary	4.442	60.542	0.000	0.000
f0 post boundary	2.654	59.139	0.010	0.040
f0 variance post boundary	2.289	59.447	0.026	0.104
prominence post boundary	2.468	68.079	0.016	0.064

nuclei.

An independent t-test grouped by IP and phrase boundary (PH) is shown in Table 1. Although significant the factors are only marginally so and only f0 pre- and post- boundary factors remained significant after bonferroni correction. If we examine the cell means in figures 3 and 4 the results are in line with previous published results. We see higher initial f0 values for after IPs, more f0 variance and more prominence caused by amplitude and duration.

If we use these factors in a discriminant analysis, we find we can categorise 72.8% of the data (71.6% with cross validation), see Table 2. Given the lack of pause data, this is in line with previous studies.

Table 2: Results of discriminant analysis using acoustic cues

Discriminant Analysis	Classification	
	PH	IP
Original	36	12
	10	23

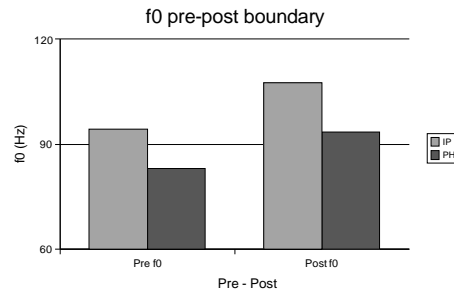


Figure 3: F0 across IP boundary and phrase boundary (PH).

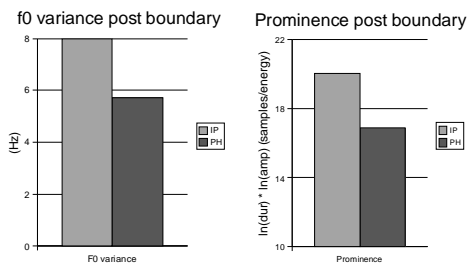


Figure 4: F0 variance and prominence -nucleus $\ln(\text{duration}) \ln(\text{amplitude}) \ln(\text{energy})$ - after IP and phrase boundary (PH).

6. Conclusion

Results show that NLPA can be used for characterising disfluency. Furthermore, that it would seem to perform as well, or better, than human subjects given the same task. Perhaps the most interesting feature of the work is that NLPA offers a non-lexical structure for dealing with timing. Using the syllable nucleus we can implicitly scale f0 contours which might allow a more structured approach to characterising intonation non-lexically. Although the prominence feature presented in this work is perhaps an over simplification of the perceptual effect of duration and amplitude, it does allow a starting point for an improved system. Similarly it would be an interesting idea to replace the f0 variance with a more perceptually based model of accentedness.

However, the success of NLPA depends largely on the autosyllabification process. Overgeneration of syllables and overestimation of syllable nuclei, for example, caused by liquids or nasals, could present a significant problem in terms of aligning f0 contours with the output. In future work we will evaluate the syllabification algorithm quantitatively against state-of-the-art autosegmentation. In addition, other acoustic features, perhaps based on spectral entropy or spectral tilt, could also be added to the system. Finally, there is a possibility that the prosodic structure produced by NLPA might be more functionally valid than one using lexical data where syllables are, in general, prescriptively assigned.

The IP analysis reinforces findings from previously published work. The results for automatic disambiguation (especially given the lack of pause information) are promising. However, in order to really test how useful these factors are for discrimination, we must also see to what extent they can tell any boundary (syllable/word) from an IP. In addition, as pointed out by Hirschberg et al [4], different speakers have different characteristics in terms of hyper-articulation. On this basis further work requires the analysis of many more subjects.

7. References

- [1] Bard, E., Lickley, R.J. & Aylett, M.P. 2001. Is Disfluency Just Difficult? *Proceedings of DISS 01, ISCA Tutorial and Research Workshop*. Edinburgh.
- [2] Dhillon, R., Bhagat, H., Carvey, H. & Shriberg, E. 2004. Meeting Recorder Project: Dialog Act Labelling Guide. *Technical Report TR-04-002*. ICSI.
- [3] Hironymous, J.L., McKelvie, D. & McInnes, F.R. 1992. Use of Acoustic Sentence Level and Lexical Stress in HSMM Speech Recognition. *ICASSP '92 Proceedings*. San Francisco. California, pp225-227.
- [4] Hirschberg, J., Litman, D. & Swerts M. 1999. Prosodic Cues to Recognition Errors. *ASRU-99*.
- [5] Howitt A.W. 2000. *Automatic Syllable Detection of Vowel Landmarks*. PhD Thesis, MIT.
- [6] Janin, A., Baron, D., Edwards, J., Ellis, D., Gelbart, D., Morgan, N., Peskin, B., Pfau, T., Shriberg, E., Stolcke, A. & Wooters, C. 2003. The ICSI Meeting Corpus. *Proceedings ICASSP-03*. Hong Kong.
- [7] Lickely, R.J. 1994. *Detecting Disfluency in Spontaneous Speech*. PhD Thesis, University of Edinburgh.
- [8] Mermelstein, P. 1975. Automatic segmentation of speech into syllabic units. *JASA*. 58(4) pp880-883.
- [9] Pannekamp A., Toepel, U., Alter, K., Hahne, A. & Friederici, A.D. 2005. Prosody-driven Sentence Processing: An Event-related Brain Potential Study. *Journal of Cognitive Neuroscience*. 17(3) pp407-421.
- [10] Tamburini, F. & Caini, C. 2005. An Automatic System for Detecting Prosodic Prominence in American English Continuous Speech. *International Journal of Speech Technology*. 8 pp33-44.