ACOUSTIC PARAMETERS OPTIMISED FOR RECOGNITION OF PHONETIC FEATURES

Anya Varnich Hansen

Center for PersonKommunikation, Aalborg University, DK-9220 Aalborg, Denmark. Tel. +45 96358671, FAX: +45 98151583, E-mail: avh@cpk.auc.dk

ABSTRACT

Speaker variability is a major problem in today's stateof-the-art speech recognition systems. Parameterisation of speech in terms of Acoustic Parameters (APs) motivated by phonetic feature theory has shown to be more robustness to speaker variability as compared to cepstral coefficients when tested on the task of broad-class recognition [1]. Also APs has been successfully applied for identification of semivowels [2,3].

The aim of the present study is to investigate the use of APs for phoneme recognition. An extended set of features is used to distinguish between all phonemes in the TIMIT database and APs related to the extended feature set are found in literature. A separability measure is calculated to investigate the importance of the suggested APs for the separation of phonemes and feature classes.

Results show that the APs that are the most important for separation of classes of phonetic features are also the most important for separation of phonemes classes. This indicates that phonemes can be recognised on the basis of phonetic features captured by the use of APs. However much work still needs to be done to understand and reliably extract all of the acoustic correlates of the phonetic features applied.

1. INTRODUCTION

Phonetic features [4] are the minimal units needed to distinguish among the most similar sounds in a language (for instance the feature voiced distinguishes the stops /b/ and /p/).

The acoustic properties of a given phonetic feature can be recognised from the speech signal by use of a carefully selected set of Acoustic Parameters (APs) [1,2,3,5]. Energy in different frequency bands measured relative to the energy in other frequency bands, the rapidity with which the spectrum changes and measurements of the relative placement of the formant frequencies are all examples of APs.

A representation of speech signals in terms of APs motivated by phonetic feature theory have shown to be more robust to inter-speaker variability as compared to cepstral based parameters when tested on the task of broad-class recognition [1]. The robustness to speaker variability is obtained by defining the APs in a relational manner across time and frequency as described above. A parameterisation of the speech signal that reduces the inter-speaker variability may result in Automatic Speech Recognition (ASR) systems that are more robust to speakers not included in the training data. An additional advantage may be minor demands to the amount of training data due to less variability.

Describing the speech signal in terms of APs motivated by phonetic feature theory also allows a directly incorporation of speech knowledge in the recognition and it serves as a tool for understanding acoustic phonetics and contextual variability [2].

The present study combines and elaborates on principles described in other feature based approaches. The fact that phonetic features manifest themselves in a set of APs is known from acoustic phonetics [5], and the principle of identifying features on the basis of APs is successfully applied in a broad-class recognition task in [1] and for identification of semivowels in [2,3]. The aim of the work presented in this paper is to be able to recognise phonemes on the basis of APs.

A set of phonetic features on the basis of which the phonemes in the TIMIT database are uniquely described and a set of APs related to the presence of the different phonetic features are found. A separability measure is applied to evaluate the importance of the suggested APs for the distinction between phonemes and phonetic feature classes.

2. MATERIALS AND METHODS

2.1 Phonetic feature set

A set of phonetic features that distinguish between all phonemes in the TIMIT database was designed on the basis of the theory described in [4]. The phonetic feature set contains the features: sonorant, syllabic, consonantal, high, low, front, back, anterior, coronal, round, voice, continuant, nasal and strident. Each of the phonemes in the TIMIT database are uniquely described in terms of these features. As an example the phoneme /aa/ is described as: [+sonorant, +syllabic, -consonantal, -high, +low, -front, +back, 0anterior, 0coronal, -round, 0voice, 0continuant, 0nasal, 0strident]. The phonetic features having the value 0 are only relevant for distinction of consonants, i.e. they are not relevant for vowels.

2.2 Acoustic parameters

Table 1 shows the applied phonetic features and their related APs found in literature [1,2,5,6,7,8,9]. Several other APs than the ones mentioned in Table 1 can be found. The *consonantal* feature can be captured from measuring the rapidity with which the spectrum changes. Rapid spectrum changes distinguish consonants from vowels and glides where the spectrum changes much more slowly [5]. In [2,3] the first difference of adjacent spectra based on a bank of 40 linear critical band filters are used as AP. In [1] APs for identifying non-*continuant* sounds are given and in [7] an auto-correlation function is suggested for the detection of *strident* sounds, see also [9] for details about this feature.

In the present study focus has initially been set on APs based on energies (E) in different frequency bands, voicing probability, the relative location of the first three formants (F1, F2, F3) and the formants bandwidth (bw1, bw2, bw3, bw4), to investigate the separability of this simple set of APs before more complex measures² are derived.

Feature	Acoustic parameter	Property /				
		reference				
sonorant	E(100-400Hz)	high [1]				
	E(0-2kHz)-E(2-8kHz)	high [1] ¹				
	E(0-300Hz)-E(3.7-7kHz)	high [2]				
	voicing probability (VP)	high [1]				
syllabic	E(600-2800Hz)	high $[1,2]^2$				
	E(2000-3000Hz)	high [1,2] ²				
	E(0-8kHz)	high [5,7]				
consonantal	E(0-8kHz)	low [6]				
	bandwidths:	broad [7]				
	bw1, bw2, bw3, bw4					
high	E(0-500kHz)-E(500-1000)	high [5]				
	F1-F0	low [2]				
back	E(800-1900kHz) - E(2-3kHz)	high [5]				
	F2-F1	low [2]				
low	E(0-500kHz)-E(500-1000)	low [5]				
front	E(800-1900kHz) - E(2-3kHz)	low [5]				
	F2-F1	high [2]				
round	no APs derived					
anterior	no APs derived					
coronal	no APs derived					
voice	E(0-400Hz)	high [6,7]				
continuant	no APs derived					
nasal	abs(250Hz-F1),	low [7,8]				
	abs(2500Hz-F2),					
	abs(3250Hz-F3)					
strident	E(6-8kHz)	high [6,9]				

Table 1 The applied phonetic features and their related APs found in literature. All energy measures (E) are calculated relative to the maximum energy across the utterance.

For some features no APs have been extracted yet as the descriptions in literature must first be transformed to measurable APs. This is the case for the feature *round*: "as sounds become more rounded, the frequencies of the higher formants decrease. But the situation is complicated in that the effect is greater in the third formant for front vowels, and in the second formant for back vowels" [8] pp.196. In this case it would be optimal to normalise with respect to a front/back measure.

For the features *anterior* and *coronal* no description of how they manifest themselves in the acoustic signal has been found in literature by now.

2.3 Extraction of acoustic parameters

Formants and bandwidth were extracted from the acoustic speech signal by use of the 'formant' tool in the Entropic software package [10]. The energy in different frequency bands was calculated by use of the Matlab signal processing tool 'specgram' [11] that computes the windowed discrete-time Fourier transform of a signal using a sliding window. A 10 ms window and a 5 ms frame rate was applied. The logarithm was taken to transform the energy measures to

dB and all energies were normalised with respect to the maximum energy across the.

2.4 The V-ratio - a separability measure

Different methods [13,14,15] can be used to determine which of the 20 APs that best distinguish between feature and phoneme classes. In the present study a seperability measure: the V-ratio that are closely related to Fishers linear discriminant function [13] has been chosen.

The V-ratio is a measure for how well an AP i separates the different classes *C*. The higher V-value, the better separation between classes.

$$V_{i} = \frac{\text{deviation of the means (over all classes)}}{\text{mean of the deviation (within each class)}} = \frac{\sigma \mu_{i}}{\mu_{\sigma_{i}}}$$

$$V_{i} = \frac{\sqrt{\frac{1}{C} \sum_{k=1}^{C} (\mu_{i,k} - \mu_{i})^{2}}}{\frac{1}{C} \sum_{k=1}^{C} \sigma_{i,k}}$$

As opposed to the methods described in [13,14] this method of selecting the most discriminative APs has the advantage that it does not require training of classifiers for the features and phonemes investigated. The drawback is that the selection of APs by use of the V-ratios does not take into account the recognition rate of the classifier.

To compare the V-ratio method to the method in [13] an experiment was performed selecting the 24 most important coefficients of a set of 12 cepstral coefficients including energy and the delta and delta-delta coefficients (in total 39 coefficients) by use of the V- ratio method. 19 of the 24 coefficients selected was the same as those selected by the method described in [13].

The ranking experiment was performed on the SISX training sentences in the TIMIT database [12]. The selection of coefficients performed in [13] was also based on the TIMIT database, but the applied dataset was not given.

2.5 Database

Data from the TIMIT database [12] sampled at 16kHz is applied for evaluation of the importance of the different APs through calculation of V-ratios. TIMIT contains 5 hours of speech and is divided into a test and a training corpus. The 3696 SISX training sentences are applied in the ranking experiments described below.

3. EXPERIMENTS

3.1 Ranking of APs according to their V-ratio

The APs are ranked according to their importance for distinction between:

- 1. +feature, -feature and Ofeature classes
- 2. phonemes (each phoneme constitutes a class)

The first ranking experiment is performed for each of the 14 features. For the features sonorant, syllabic, consonantal, high, back and low there are two classes: the

¹ In [1] the AP is not normalised with the maximum energy across the utterance.

 $^{^2}$ In [1,2] to capture the syllabic feature an energy peak is detected and measured relative to surrounding energy deeps in the neighbouring consonant - for details see [3]. E.i. in the present study a simplified measure has been applied.

Feature →	son	syl	cons	high	back	low	fro	ant	cor	rou	voi	cont	nas	str	pho-
AP		-													neme
(designed to capture the															
feature in brackets)															
\downarrow															
1 p-voi	2.000	1.424	1.544	0.831	0.955	1.245	1.495	0.329	0.248	1.333	0.644	0.461	0.418	0.669	1.499
(sonorant)	1	3	1	6	5	5	3	1	4	3	1	2	6	4	1
2 bw1	0.316	0.246	0.249	0.097	0.191	0.285	0.349	0.124	0.191	0.319	0.242	0.197	0.234	0.376	0.560
(consonantal)	13	14	13	18	14	15	13	15	5	11	10	13	12	8	14
3 bw2	0.264	0.208	0.236	0.172	0.136	0.180	0.315	0.188	0.133	0.245	0.176	0.186	0.179	0.266	0.306
(consonantal)	15	15	14	15	16	17	14	11	10	15	13	15	15	11	17
4 bw3	0.097	0.099	0.104	0.097	0.078	0.080	0.103	0.038	0.022	0.105	0.040	0.044	0.035	0.045	0.150
(consonantal)	18	19	18	19	19	19	20	19	19	18	20	20	20	20	19
5 bw4	0.070	0.077	0.089	0.084	0.058	0.058	0.106	0.022	0.019	0.081	0.053	0.049	0.042	0.058	0.131
(consonantal)	19	20	19	20	20	20	19	20	20	20	19	19	19	18	20
6 F1-F0	0.610	0.519	0.551	0.436	0.426	0.478	0.640	0.197	0.148	0.606	0.371	0.260	0.356	0.457	0.603
(high/low)	11	10	10	8	9	11	10	10	9	10	8	9	9	6	13
7 F2-F1	0.023	0.112	0.050	0.152	0.119	0.264	0.115	0.105	0.116	0.085	0.102	0.102	0.254	0.057	0.383
(back/front)	20	18	20	17	18	16	18	18	13	19	18	18	10	19	16
8 abs(250-F1)	0.631	0.522	0.558	0.426	0.413	0.476	0.664	0.206	0.151	0.627	0.405	0.265	0.382	0.459	0.629
(nasal)	10	9	9	9	11	12	9	9	8	9	6	8	8	5	11
9 abs(2500-F2)	0.296	0.322	0.275	0.243	0.294	0.450	0.356	0.111	0.130	0.258	0.152	0.189	0.242	0.138	0.552
(nasal)	14	12	12	12	13	13	12	17	11	12	14	14	11	17	15
10 abs(3250-F3)	0.154	0.157	0.195	0.161	0.121	0.118	0.225	0.146	0.113	0.249	0.140	0.141	0.200	0.158	0.261
(nasal)	16	17	16	16	17	18	15	14	14	14	15	17	13	16	18
11 E(100-400)	1.620	1.446	1.473	1.035	1.086	1.370	1.518	0.212	0.107	1.425	0.463	0.209	0.448	0.343	1.353
(sonorant)	2	1	2	2	2	2	1	8	15	1	5	12	3	10	2
12 E(0-2k)-E(2-8k)	0.836	0.607	0.589	0.303	0.511	0.702	0.800	0.254	0.309	0.735	0.578	0.377	0.537	0.882	1.244
(sonorant)	7	8	8	13	8	7	7	3	3	8	3	4	2	3	4
13 E(0-300)-E(3.7-7k)	0.791	0.507	0.533	0.268	0.418	0.502	0.727	0.244	0.332	0.743	0.634	0.445	0.725	0.934	1.197
(sonorant)	8	11	11	11	10	10	8	4	2	7	2	3	1	1	6
14 E(600-2800)	1.273	1.320	1.328	0.979	1.013	1.340	1.261	0.232	0.048	1.127	0.229	0.324	0.186	0.200	1.206
(syllabic)	5	5	5	4	4	4	4	5	17	4	12	6	14	12	5
15 E(2-3k)	1.009	1.098	1.131	0.959	0.926	1.050	0.917	0.224	0.041	0.816	0.105	0.258	0.096	0.258	0.896
(syllabic)	6	6	6	5	6	6	6	6	18	6	17	10	17	12	8
16 E(0-8k)	1.307	1.436	1.421	1.145	1.151	1.445	1.145	0.181	0.098	1.058	0.107	0.371	0.064	0.247	1.155
(consonantal)	4	2	4	1	1	1	5	12	16	5	16	5	18	13	7
17 E(0-500)-E(0.5-1k)	0.106	0.201	0.162	0.174	0.147	0.357	0.194	0.118	0.155	0.164	0.237	0.300	0.447	0.224	0.612
(high/low)	17	16	17	14	15	14	16	16	6	16	11	7	4	14	12
18 E(0.8-1.9)-E(2-3k)	0.376	0.306	0.229	0.271	0.298	0.535	0.401	0.148	0.153	0.257	0.300	0.156	0.136	0.408	0.860
(front/back)	12	12	15	10	12	9	11	13	7	13	9	16	16	7	10
19 E(0-400)	1.612	1.423	1.456	1.012	1.066	1.343	1.517	0.216	0.116	1.420	0.471	0.222	0.446	0.371	1.344
(voiced)	3	4	3	3	3	3	2	7	12	2	4	11	5	9	3
20 E(6-8k)	0.672	0.652	0.639	0.621	0.664	0.673	0.170	0.259	0.369	0.118	0.372	0.570	0.395	0.899	0.871
(strident)	9	7	7	7	7	8	17	2	1	17	7	1	7	2	9

Table 2 For each of the APs suggested in literature a separability measure (the V-ratio) is calculated, to determine the importance of the AP in the separation of for instance the +sonorant class of phonemes from the -sonorant class of phonemes (see the third column). For each of the features V-ratios are given together with the resulting ranking of the APs. In the last column each phoneme constitute a class and the V-ratios describe the importance of a given AP for the separation of all phonemes. The numbers given in the left column are used when referring to the different APs in the discussion section, below the AP number is given the name of the feature it was designed to capture.

class of phonemes in which the feature is present (denoted +feature) and a class of phonemes where the feature is not present (denoted -feature). For the remaining features the phonemes for which the feature is not relevant constitutes a third class denoted the 0feature class. In the second ranking experiment each of the phonemes constitutes a phoneme class.

The first ranking experiment will show whether APs expected on the basis of acoustic phonetic knowledge to be important for the caption of a certain features are important according to V-ratios calculated on a large set of data. The second ranking will indicate which features to use when training phoneme models for speech recognition.

4. RESULTS

Results are given in Table 2. For each feature the Vratios are given for each of the APs together with the ranking given as numbers from 1 to 20 where 1 is the most important AP for the given feature.

5. DISCUSSION

Comparing the ranking of APs for the features and phonemes it is seen, that it is the same APs (see grey shaded areas in Table 2) that are important for recognition of features and phonemes. This indicates that phonemes can be recognised on the basis of phonetic features captured in the acoustic speech signal by the use of APs. For all features relevant for the distinction of vowels e.i. son, syl, cons, high, back, low, fro, rou the APs 1, 11, 14, 16 and 19 have high V-values. All these APs are highly related to the detection of the feature *sonorant*, which split phonemes into the two major classes: sonorants and obstruents. The presence of one phonetic feature can give information on the presence of absence of others. Thus if a phoneme is -sonorant, it is known that it is also *-syllabic*, *+consonantal* and *-low* therefore the same set of APs are important for these features, for the distinction into major classes, however some of the APs with lower V-values may be responsible for the fine phonetic distinctions.

The voicing probability (AP1), AP 11 and 19 capturing the strong low frequency energy resulting from voicing are important APs for the *sonorant* feature. The presence of formants in the mid frequencies (AP14) and a high overall energy (AP16) are also resulting in high V-values for these APs. AP 12 and 13 were designed to capture the sonorant feature as having a strong low-frequency energy as compared to the energy in the higher frequencies. However these APs do not have as high V-ratios as the ones first mentioned. Some of the APs such as AP1, 11 and 19 are highly correlated, therefore a constraint should be set on the maximum allowable correlation among APs when they are selected for training of for instance a HMM classifier.

AP 14 and 15 that was designed to capture *syllabic* peaks in the mid-frequency regions get relatively high V-ratios and thus the numbers 5 and 6 in the ranking. Even higher V-ratios might be obtained with a more advanced APs^2 . Though written in [7] that "phonemes possessing the *consonantal* feature are acoustically characterised by a broadening, reduction and fusion of formants and formant regions due to zeros, high damping or transient variations of formant frequencies" the bandwidths (AP2, 3, 4, 5) do not seem to be relevant APs.

Different values of the same AP can characterise different features, thus AP16 is high for sonorant sounds and low for consonantal sounds that are characterised by having a low total energy.

Of the APs designed to capture the features *high/low* AP6 is superior to AP17. While for the features *front/back* AP18 have higher V-ratios than AP7.

No APs was designed to capture the features *anterior*, *coronal*, *round* and *continuant*, so it is not surprising that the V-ratios for these features in general are low. The feature *round* is an exception. This feature is only relevant for vowels, therefore the APs that captures the characteristics of vowels (e.i. the +sonorant, +syllabic, - consonantal features) are given a high priority. The high V-ratios for these APs may express that vowels are well separated from other phonemes, thus the high V-ratios do not ensure, that rounded and unrounded vowels are well separated.

Of the APs 8, 9 and 10 suggested for capturing the feature *nasal* AP8 shows to be the most important. The very low first formant centred at about 250Hz, that characterise nasal sounds [8] is also captured by AP1, 11, 12, 13, 17 and 19. AP20 designed to capture the strident feature show to be highly relevant and gets a priority as the second most important feature.

In general the features that are not relevant for vowels (ant, cor, voi, cont, nas, str) have low V-values indicating that more efforts must be put in finding APs that better capture these features that perform the fine phonetic distinctions. As opposed to the cepstral representation of speech it is possible to put special efforts in capturing 'difficult feature'.

6. CONCLUSION

Results show that the APs that are the most important for separation of classes of phonetic features are also the most important for separation of phonemes classes. This indicates that phonemes can be recognised on the basis of phonetic features captured by the use of APs. However much work still needs to be done to understand and reliably extract all of the acoustic correlates of the phonetic features applied.

7. REFERENCES

- Bitar, N. and Espy-Wilson, C. (1996) "Knowledge-based parameters for HMM speech recognition", ICASSP'96, pp 29-32.
- [2] Espy-Wilson, C. Y. (1994) "A feature-based semivowel recognition system", J. Acoust. Soc. Am. 96 (1), July 1994, pp.65-72.
- [3] Espy-Wilson, C. Y. (1992) "Acoustic measures for linguistic features distinguishing the semivowels /wjrl/ in American English", J. Acoust. Soc. Am. 92 (2), Pt.1, August 1992, pp.736-757.
- [4] Chomsky, N. and Halle, M. (1968) "The Sound Pattern of English", Harper and Row, New York.
- [5] Stevens, K. (1980) "Acoustic correlates of some phonetic categories", J. Acoust. Soc. Am. 68 (3), pp. 836-842.
- [6] Clark, J. and Yallop, C. (1991) "An introduction to phonetics and phonology", Basil Blackwell, Inc. Cambridge, Massachusetts, USA.
- [7] Jakobson, R., Fant, G., Halle, M. (1951) "Preliminaries to speech analysis", The MIT Press.
- [8] Ladefoged, P. (1993) "A course in phonetics", Harcourt Brace Jovanovich College Publishers.
- [9] Utman, J.A. and Blumstein, S.E. (1994) "The Influence of Language on the Acoustic Properties of Phonetic Features: A Study of the Feature [Strident] in Ewe and English", Phonetica 1994; 51 pp.221-238.
- [10]"Waves+ Manual" (1993) version5.0, Entropic Research laboratory, Inc.
- [11]Krauss, P. (1993) "Signal Processing TOOLBOX for use with Matlab", The Math Works, Inc.
- [12]Garofolo, J. et. al. (1993) "DARPA TIMIT Acoustic-Phonetic Continuous Speech Corpus CD-ROM" U. S. Department of Commerce, National Institute of Standards and Technology, USA
- [13]Bocchieri, E.L. and Wilpon, J.G.(1993) "Discriminative feature selection for speech recognition", Computer Speech and Language, 1993, 7, pp.229-246.
- [14]Phillilps, M. and Zue, V. (1992) "Automatic Discovery of Acoustic Measurements for Phonetic classification", ICSLP'92, vol. 1, pp 795-798.
- [15]Duda, R. and Hart, P. (1973) "Pattern Classification and Scene Analysis", John Wiley and Sons, Inc pp 114-121.