

Some Investigations of the Decision Trees from the ASR, TTS and Applied Phonetics Point of View

Stefan Grochowski

Poznan University of Technology, Poznan, Poland

E-mail: stefan.grochowski@cs.put.poznan.pl

ABSTRACT

In our HMM's based ASR system we have used the models of Polish triphones. To deal with triphones for which there are no examples in the training data we introduced decision tree-based clustering. In the paper we tried to analyze the clusters with the identical triphones and the questions set forming the decision trees obtained in our experiments with the ASR system. The results of our investigations of the equivalent contexts leading to the acoustically similar triphones (they share common clusters for all states) should be useful from the point of view of speech recognition, text to speech systems and applied phonetics.

1. INTRODUCTION

The decision tree-based states clustering mechanism offers a solution to the unseen triphone problem. A phonetic decision tree is a tree in which a binary phonetic question is attached to each node. Initially all states are placed in a single cluster at the root of the tree. The question at each node is chosen in order to maximize the likelihood of the training data. Splitting any node into two will increase the log likelihood. The increase obtained when each possible question is used can be calculated and the "best" question selected which gives the biggest improvement [3]. The process of node splitting is stopped if the log likelihood gain becomes lesser than the predefined threshold. The states in each cluster are then tied. After tying procedure we can obtain many identical triphones, with identical all appropriate states. In this paper we will study which contexts are equivalent, i.e. result in the same states. It will help us to compare our theoretic phonetic knowledge with the real data. We will also investigate the sequences of questions in the decision trees since their succession is connected with their weights. In the section 2 we will present the database for Polish used in our experiments and in the section 3 the details about the acoustic modeling will be presented.

2. DATABASE FOR POLISH

The investigation has been made on the base of the first database for Polish CORPORA [2]. It was created to comprise the relatively large number of Polish diphones.

The final number of diphones was equal to 179.753 and the number of different diphones we have found in CORPORA was equal to 1271. Note that 37 Polish phonemes implies about 1440 theoretically possible diphones. The database is distributed with the ASCII files in HTK [3] format comprising the phonemic transcriptions along with the timing data.

3. ACOUSTIC MODELING OF POLISH

We have chosen the context dependent phonemes (triphones) as the phonetic units and Continuous Density Hidden Markov Models (CDHMM) as the models. For training models of 37 Polish phonemes we have used 5110 utterances of CORPORA (14 male speakers x 365 utterances). We found 2953 different triphones which give 8859 states (3 states per model). The biggest problem in building context dependent HMM's is data insufficiency. Considering the size of our database we can expect that not all triphones will be reliably estimated. The solution consists in tying the states within triphone sets in order to share data and thus be able to make robust parameter estimates. The choice of which states to tie requires a little subtlety since the performance of the recognizer depends crucially on how accurately the state output distributions capture the statistics of the speech data. In our case we have chosen the decision tree based approach [3] to tie the similar states. This approach helps to solve the problem of the unseen triphones, i.e. it allows triphone models which have no training data to be synthesized. The phonetic decision trees are built using a top-down sequential optimization process. Initially all states are placed in a single cluster at the root of the tree. The question is then found which gives the best split of the root node. This process is repeated until the increase in log likelihood falls below the specified threshold. As a final stage, the decrease in log likelihood is calculated for merging terminal nodes with different parents. Any pair of nodes for which this decrease is less than the threshold used to stop splitting are then merged.

In the section 4 we will investigate the sequence of questions during the decision tree construction.

After tying procedure we can expect that different triphones from the same monophone can share the same corresponding states. In the section 5 we will investigate the identical triphones. It will be interesting to know

which contexts are equivalent (only for our database of course) and if this knowledge corresponds to the acoustic phonetic theory.

4. ANALYSIS OF THE QUESTIONS IN THE DECISION TREES

In three states models, in very coarse approximation, we assume that the left state is left context dependent,

phoneme	questions about context (in %)				
	state 1	state 2		state 3	
	left	left	right	right	
1	<i>i</i>	100	34,5	65,5	70,8
2	<i>i</i>	100	65	35	73,7
3	<i>e</i>	85,2	50,9	49,1	75
4	<i>a</i>	84,6	46,9	53,1	90
5	<i>o</i>	88	51,9	48,1	92,9
6	<i>u</i>	86,7	53,3	46,7	80
7	<i>õ</i>	90	77,8	22,2	36,4
8	<i>ẽ</i>	50	28,6	71,4	75
9	<i>w</i>	60	71,4	28,6	58,3
10	<i>j</i>	88,2	75	25	45,4
11	<i>l</i>	60	52,6	47,4	50
12	<i>t</i>	87,5	45,5	54,5	50
13	<i>p</i>	64,3	80	20	70
14	<i>k</i>	70	46,2	53,8	54,5
15	<i>d</i>	100	80	20	88,9
16	<i>b</i>	53,8	14,3	85,7	100
17	<i>g</i>	63,6	50	50	100
18	<i>ts</i>	87,5	42,9	57,1	50
19	<i>tf</i>	100	50	50	75
20	<i>tç</i>	88,9	20	80	50
21	<i>dz</i>	0	0		0
22	<i>dž</i>	100	0	100	100
23	<i>dž</i>	100	0	100	100
24	<i>m</i>	55,6	52,9	47,1	66,7
25	<i>n</i>	87,5	40	60	72,7
26	<i>ɲ</i>	100	71,4	28,6	66,7
27	<i>ŋ</i>	100	100	0	100
28	<i>s</i>	90	42,1	57,9	100
29	<i>ʃ</i>	71,4	55,6	44,4	71,4
30	<i>ç</i>	100	66,7	33,3	85,7
31	<i>f</i>	80	83,3	16,7	75
32	<i>x</i>	83,3	50	50	75
33	<i>r</i>	72,7	50	50	58,3
34	<i>w</i>	50	33,3	66,7	85,7
35	<i>z</i>	87,5	40	60	87,5
36	<i>ž</i>	100	33,3	66,7	83,3
37	<i>ž</i>	100	100	0	100
		83,0	51,5	48,5	75,4

Table 1: Results of counting the questions for all states about the left or right context

the middle one is context independent and the right one is

right context dependent. To verify this opinion in the Table 1 we present the results of counting the questions about the appropriate context. We can see, for example, for the vowel *i* that the left state always (for our database) depends only on the left context (100%) whereas the right state of the model depends not only on the right context (in 70,8%) but also on the left one. Nevertheless, studying the Table 1 we can state that in most cases the boundary states depend also on the opposite contexts.

We have found some very interesting exceptions. Splitting the left state of the consonant *v* on the first step, before asking about the left context, we should ask if the right context is *i*. We have found the same situation splitting the right state of the fricative *ʃ*. The main question was about the left context (is it *u*?). A similar case was for fricative *f*. Splitting the right state we should ask if the left context is *a*. To complete the exceptions we should mention the right state of the plosive *k* (depends mainly on the *e* on the left side) and *w* - where the left state depends mainly on the *a* on the right side.

When interpreting the above results we have been conscious that there was a dependence between the implementation of the particular question for the particular phone and the number of data observations for that phone. If the database used for the design of the phone set has phone frequencies which are very different from the actual recognition system environment, the resulting decision tree and the designed phone set may be different.

	Number of occurrences	Number of different triphones		Number of occurrences	Number of different triphones
<i>a</i>	5857	222	<i>b</i>	966	65
<i>e</i>	4043	220	<i>f</i>	942	50
<i>o</i>	3368	192	<i>ʃ</i>	862	61
<i>n</i>	2561	114	<i>g</i>	789	56
<i>j</i>	2205	98	<i>ts</i>	779	57
<i>i</i>	2193	127	<i>z</i>	772	67
<i>r</i>	2070	103	<i>ž</i>	759	67
<i>m</i>	1914	108	<i>x</i>	740	53
<i>t</i>	1828	80	<i>õ</i>	715	59
<i>l</i>	1772	86	<i>ç</i>	661	66
<i>i</i>	1722	105	<i>tç</i>	588	52
<i>u</i>	1689	97	<i>tf</i>	508	43
<i>k</i>	1514	79	<i>z</i>	450	43
<i>w</i>	1447	83	<i>dž</i>	383	37
<i>s</i>	1319	86	<i>ẽ</i>	315	37
<i>w</i>	1235	73	<i>dž</i>	247	23
<i>d</i>	1167	85	<i>dz</i>	160	18
<i>p</i>	986	59	<i>ŋ</i>	108	16
<i>ɲ</i>	967	64	Σ	52418	2953

Table 2: Statistics of the dataset used in the experiments

In the table 2 we present the statistics of our dataset; the number of occurrences of the particular phonemes along with the number of different triphones for that phoneme.

In the table 3 we present the most important questions about the left and right context for all phonemes.

		left state	right state
1	<i>i</i>	<i>j,n,m,l,dz,v</i>	<i>w,m,n,k,o,r,z</i>
2	<i>ɪ</i>	<i>m,w,n,r,ʒ,b</i>	<i>w,,n,m,j,s,k,b</i>
3	<i>e</i>	<i>j,j,l,dz,tɕ,ʒ</i>	<i>j,n,m,j,w,o,r</i>
4	<i>a</i>	<i>j,n,w,m,l,r,d</i>	<i>n,j,r,w,l,m</i>
5	<i>o</i>	<i>r,n,j,e,d,w</i>	<i>j,n,l,j,m,k,s,r</i>
6	<i>u</i>	<i>j,m,l,e,r,g</i>	<i>j,ʃ,l,r,ʒ,ts</i>
7	<i>õ</i>	<i>j,n,d,s,ʒ</i>	<i>e,u,ʒ,i,s,dj</i>
8	<i>ẽ</i>	<i>ɕ,s</i>	<i>s,ʃ,ɕ,z</i>
9	<i>w</i>	<i>a,a,e,i,i</i>	<i>a,o,i,u,g,z,u</i>
10	<i>j</i>	<i>a,e,u,m,o,r</i>	<i>a,o,u,o,õ,jn</i>
11	<i>l</i>	<i>a,o,e,u,i,i</i>	<i>e,a,i,u,j,o,a,o</i>
12	<i>t</i>	<i>a,o,s,r,e,n</i>	<i>a,o,e,i,s,n,r,e</i>
13	<i>p</i>	<i>a,o,u</i>	<i>a,o,i,j</i>
14	<i>k</i>	<i>i,o,e,a,i</i>	<i>a,e,o,i,i</i>
15	<i>d</i>	<i>o,a,i,e</i>	<i>a,o,e,õ,j</i>
16	<i>b</i>	<i>o,a,u,r,e,e</i>	<i>a,i,e,o,e,r</i>
17	<i>g</i>	<i>a,o,i,w,e,u</i>	<i>u,o,a,e,i</i>
18	<i>ts</i>	<i>a,i,o,u,a</i>	<i>j,e,i,a</i>
19	<i>tʃ</i>	<i>e,i,i,a</i>	<i>i,e,a,i</i>
20	<i>tɕ</i>	<i>ɕ,a,o,a</i>	<i>i,e,a,i,jn,a</i>
21	<i>dz</i>	-	-
22	<i>dʒ</i>	<i>o,a,i</i>	<i>i,e</i>
23	<i>dʒ̥</i>	<i>e</i>	<i>o,e</i>
24	<i>m</i>	<i>e,a,o,i,i,a,i,u</i>	<i>a,i,j,u,i</i>
25	<i>n</i>	<i>a,o,e,i,i,n</i>	<i>a,t,o,õ,e,n</i>
26	<i>ɲ</i>	<i>o,e,a</i>	<i>e,i,tɕ,j,e,o</i>
27	<i>ŋ</i>	<i>o</i>	<i>g</i>
28	<i>s</i>	<i>a,o,e,i</i>	<i>w,t,k,e</i>
29	<i>ʃ</i>	<i>u,a,e,i,i,i</i>	<i>u,i,a,e,k,a</i>
30	<i>ɕ</i>	<i>i,a,o,e,m</i>	<i>tɕ,i,e,ẽ</i>
31	<i>f</i>	<i>a,e,o,u,e</i>	<i>a,a,j,i,p,e</i>
32	<i>x</i>	<i>e,a,iu,o,u</i>	<i>a,o,u,w,e,u</i>
33	<i>r</i>	<i>a,u,e,o,i,u</i>	<i>t,o,a,j,u</i>
34	<i>v</i>	<i>i,e,o,a,j,d,u,b</i>	<i>a,i,o,j,u,o</i>
35	<i>z</i>	<i>i,a,i,o</i>	<i>a,j,e,i,o</i>
36	<i>ʒ</i>	<i>u,a,o</i>	<i>i,e,u,o</i>
37	<i>ʒ̥</i>	<i>a,o,d</i>	<i>i,e,a,o</i>

Table 3: The most important questions about the left and right context for all phonemes

The very general conclusion is that for vowels the most important questions were about nasals, liquids and glides whereas for consonants the most important questions

were about vowels. Letters in bold in the Table 3 mean that an appropriate state depends on the opposite context - see the beginning of the section.

5. EQUIVALENT CONTEXTS

After tying procedure for almost all phonemes we have obtained some clusters of triphones. In every cluster all triphones have identical appropriate states. In the table 4 the total number of different triphones, the number of clusters, the number of triphones outside clusters and the number of triphones per cluster are presented.

	different triphones	number of clusters	outside clusters	triphones per cluster
<i>a</i>	222	10	87%	2,8
<i>õ</i>	59	13	12%	4
<i>b</i>	65	11	28%	4,3
<i>ts</i>	57	8	28%	5,1
<i>tʃ</i>	43	0	100%	
<i>tɕ</i>	52	11	10%	4,3
<i>d</i>	85	11	40%	4,6
<i>e</i>	220	19	77%	2,8
<i>ẽ</i>	37	6	11%	5,5
<i>f</i>	50	10	8%	4,6
<i>g</i>	56	8	37%	4,4
<i>x</i>	53	10	11%	4,7
<i>i</i>	127	20	42%	3,7
<i>j</i>	98	10	37%	6,2
<i>k</i>	79	11	39%	4,4
<i>l</i>	86	9	35%	6,2
<i>w</i>	83	13	13%	5,5
<i>m</i>	108	13	21%	6,5
<i>n</i>	114	13	27%	6,4
<i>ɲ</i>	64	6	16%	9
<i>o</i>	192	19	78%	2,3
<i>p</i>	59	16	17%	3,1
<i>r</i>	103	13	18%	6,5
<i>ʒ</i>	67	10	19%	5,4
<i>s</i>	86	12	38%	4,4
<i>ʃ</i>	61	10	20%	4,9
<i>ɕ</i>	66	13	14%	4,4
<i>t</i>	80	13	36%	3,9
<i>u</i>	97	19	43%	2,9
<i>w</i>	73	10	14%	6,3
<i>ɪ</i>	105	19	45%	3,1
<i>z</i>	67	14	18%	3,9
<i>ʒ</i>	43	6	12%	6,3
<i>dz</i>	18	0	100%	
<i>dʒ̥</i>	23	0	100%	
<i>dʒ</i>	37	0	100%	

Table 4: The total number of different triphones, the number of clusters, the number of triphones outside clusters and the number of triphones per cluster

For four consonants: $tʃ$, $dʒ$, $dʒ$, $dʒ$ we have observed no clusters. It means that they are very context dependent. Almost the same can be said about three vowels: a , e , o (about 80 % outside clusters). On the opposite side we have diphthongs: $õ$, $ẽ$ and fricatives f , x .

It is interesting to study the lists of the identical triphones and triphones placed in different clusters. For example, the set $\{i+b, x-i+b, ts-i+dʒ, i+dʒ, i+dʒ, i+e, i+f, i+p, i+v, j-i, t-i+dʒ, r-i, n-i+e, a-i+dʒ, i-i+p\}$ forms one cluster, whereas triphones $i+m$ and $i+n$ form different clusters. The notation $h-i+b$ describes the phoneme i with the left context x and right context b [3]. We should be careful since our tying procedure consists in tying the states in each leaf node in the decision tree. It means that the triphones in one cluster are the same but at the same time the triphones in different clusters could be the same. For example, $d-a+m$ and $d-a+n$, $n-a+m$ and $n-a+n$ are not tied because the questions about m and n are the subsequent questions building the decision tree and they must be in different clusters. We noticed the same situation for all vowels. Thus, one can suggest asking not about the particular phonemes but about the classes as in [1]. Nevertheless, in such a solution we are not able to investigate the succession of questions. In the question about the right state of e the nasals are on the: second (n), third (m) and fourth ($ɲ$) position whereas in the question about the right state of a the nasals m and n are separated by five other questions. Considering one cluster we can draw conclusions about the equivalent contexts. For example in the cluster $\{n-a+a, n-a+e, n-a+i, n-a+o, n-a\}$ the right vowel contexts are not distinguishable.

6. DISCUSSION OF RESULTS

Analyzing the questions in the decision trees and the clusters of triphones obtained after tying the states in leaves we can enhance our knowledge of acoustic phonetics from the statistical point of view. With this technique we can discover new dependencies, especially connected with the significance of the context. This knowledge can be useful from the ASR, TTS and applied phonetics point of view. The obtained results depend on the dataset. To confirm our observations we analyzed the appropriate decision trees for another dataset with different vocabulary. The results were very similar. It is impossible to present all of them. The results are in the form of the set of questions and the contents of the leaves in the decision trees. The goal of our paper was to present rather the used "technology" than the detailed results.

7. CONCLUSION

In the paper the investigations of the decisions trees used in our HMM based ASR system have been

presented. Two kinds of data: the succession of questions and the contents of leaves have been studied. In our opinion the results of the investigation should enhance our phonetic knowledge from the ASR, TTS and applied phonetics point of view.

ACKNOWLEDGMENT

This work was supported by KBN, grant no. 7T 11C 009 21

REFERENCES

- [1] V.Chuchupal, K.Makovkin, A.Chichagov, Accurate acoustic modelling for Russian, *Proc. SPECOM'2000*, pp.71-74, 2000.
- [2] S.Grocholewski, First Database for Spoken Polish, *Proc. Int. Conf. on Language Resource & Evaluation*, pp.1059-1062, 1998.
- [3] S.Young, et al., *HTK Book*, Cambridge University, 1997.