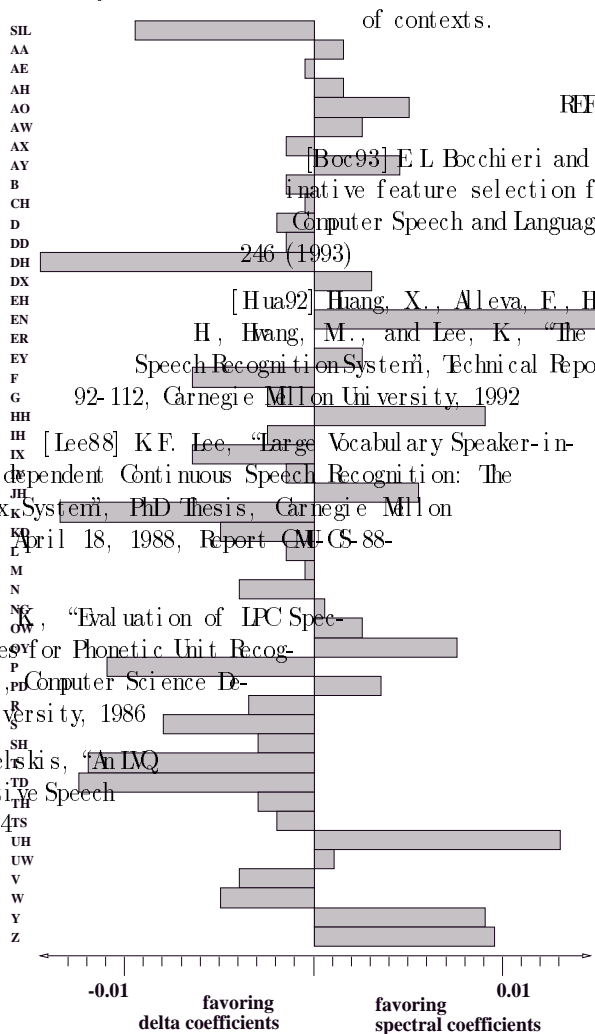


phoneme tend more to favoring delta-coefficients of features, like e.g. delta-spectral-coefficients, although one might expect that these coefficients' delta-delta-spectral-coefficients, power, acoustics are rather static and less context-dependent. Certainly, one fact that delta coefficients do model the dynamic delta-feature. Experiments with non-context-independent signal like will give also more information about the dependence of the streamweights on the different types of contexts.



REFERENCES

[Boc93] E.L. Bocchieri and J.G. Wipon, "Discriminative feature selection for speech recognition", *Computer Speech and Language*, Vol. 7, pp. 229-246 (1993)

[Hua92] H. Hwang, X. Alleva, E. Hayamizu, S. Hsu, H. Hwang, M., and Lee, K., "The SPHINX II Speech Recognition System", Technical Report CS-92-112, Carnegie Mellon University, 1992

[Lee88] K.F. Lee, "Large Vocabulary Speaker-independent Continuous Speech Recognition: The Sphinx System", PhD Thesis, Carnegie Mellon University, April 18, 1988, Report CM-CS-88-112

[Kaw86] K. Kawahara, "Evaluation of LPC Spectral Measures for Phonetic Unit Recognition", Report, Computer Science Department, Carnegie Mellon University, 1986

[Ebe84] S. Ebeliski, "An LQ Adaptive Speech Recognition System", Report CM-CS-84-44

4. FUTURE WORK

So far we have only performed experiments with streams. We believe that the proposed approach will be even more fruitful for systems with

$\alpha_i(B)$ after iteration $k + n$ to be approximately
 $\alpha_i(B)^t - n \cdot \lambda(d^*LP_i(\alpha, B)/d^*\alpha_i(B))$ (if no sig-
 mid is applied). We have found that the dif-
 ferences from iteration to iteration are in fact so
 small that this approximation is valid, which sug-
 gested a second solution to the above mentioned
 problem namely to run simply one or two itera-
 tions with a large stepsize, or alternatively to use
 a cross validation mechanism to decide what num-
 ber of iterations (i.e. what stepsize λ) is best.

3. EXPERIMENTS

We have performed experiments on the English
 Registration Task (CR) [Wo92] and
 Management Task (RM), using the
 [92] of the JANUS Speech to
 [Wi91]. The recog-
 nition probabilities for
 a 50-cluster
 probability
 300 context
 1000
 ma

path, C , did not get defined because of non-convexity of the loss function, whose domain is the set of all possible states. We are numerically minimizing the loss function, whose domain is the set of all possible states. We are numerically minimizing the loss function, whose domain is the set of all possible states.

$$LP_t(\alpha, C) := -\log P(x_i|C) = \sum_{i=1}^n c_i(t) \cdot \alpha_i(C) \quad (5)$$

For a simple two-feature system, eq. (5) results in $\alpha_j(B) \text{ updated} = \alpha_j(B) + \lambda \cdot \frac{d LP_t(\alpha, B)}{d \alpha_j(B)}$ (3)

$$\frac{d LP_t(\alpha, B)}{d \alpha_j(B)} = b_1 + \frac{\alpha_2(B) b_2}{\alpha_1(B)} - b_2 \quad (6)$$

We can easily see, in the general case the

updated system will produce a higher probability for the correct path (or for some given labels). Note that the partial derivative $\frac{\partial LP_t(\alpha, S)}{\partial \alpha_i(S)}$ would be independent of each other. Because of the above mentioned summation constraint, many gradient descent steps must be taken to reach the constraint. The probability of a state being increased must be increasing all $\alpha_i(S)$, but because we would like to diversify different features.

step $\alpha_j(B)$ is $(\alpha_1(S), \dots, \alpha_n(S))$ $(\alpha'_1(S) = \alpha_1(S) + \delta_1, \dots, \alpha'_n(S) = \alpha_n(S) + \delta_n)$ could thus be defined as $\delta_j := \epsilon$ and $\delta_i = 0$ for $i \neq j$. This step definition is the next, while all the other $\alpha_i(S)$ are unchanged. Other step

LEARNING STATE-DEPENDENT STREAM WEIGHTS FOR MULTI - CODEBOOK HMM SPEECH RECOGNITION SYSTEMS

I. Rogina, A. Wabel

University of Karlsruhe, Postfach 6980, 76128 Karlsruhe, Germany
Carnegie Mellon University, Pittsburgh, PA15213-3890, USA

1. TRAINING

ABSTRACT

which uses n information streams x_t , and for a given HMM S , the overall probability is then $P(x_t | S) = \prod_{i=1}^n P_i(x_t | S)$, where $P_i(x_t | S)$ is the probability of x_t given S using the i -th codebook. The final score is a weighted sum of the contributions of every codebook, and $\sum_{i=1}^n \alpha_i = 1$. These weights can be found empirically by the same set of weights is used for every HMM state. There is reason to believe that features which are more important for models than for others. Especially, the beginning and ending segments are more context dependent than the middle segments. In that case the probability for a stream with a spectrum that is more similar to the spectrum of the codebook coefficients should be higher than for a stream with a spectrum that is less similar to the spectrum of the codebook coefficients.