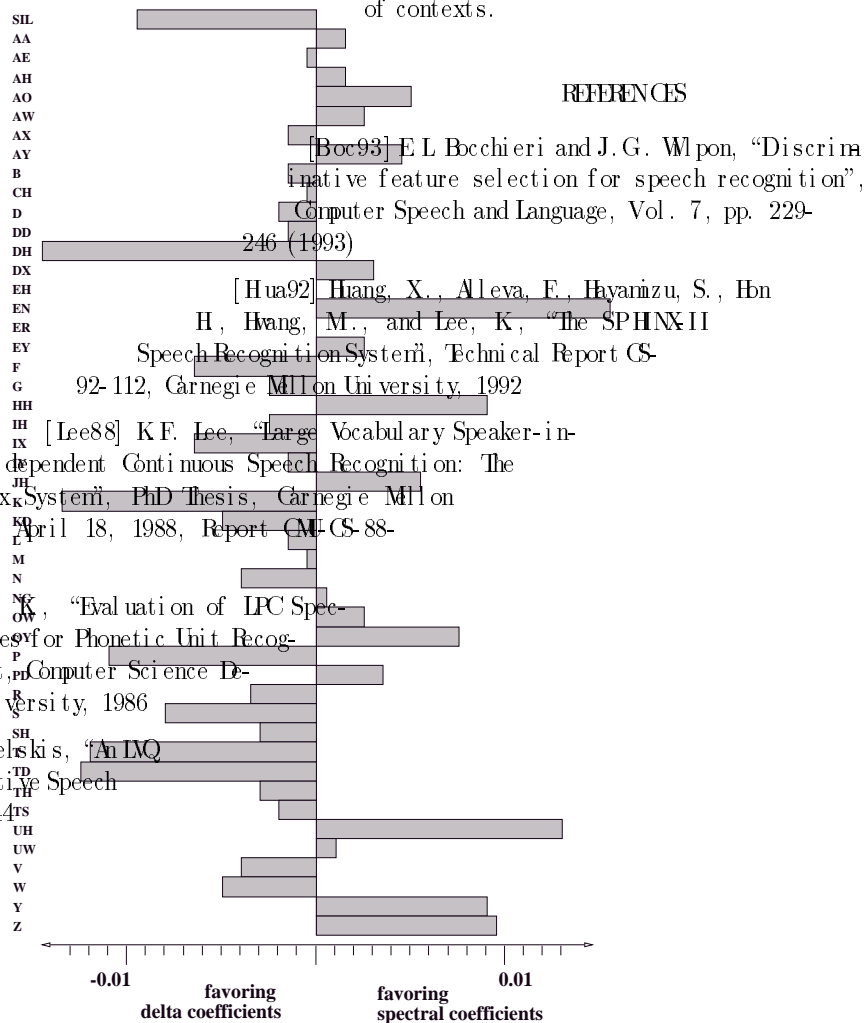


phoneme tend more to favoring delta-coefficients of features, like e.g. delta-spectral-coefficients, although one might expect that these coefficients' delta-delta-spectral-coefficients, power, acoustics are rather static and less context-dependent. Certainly, one fact that delta coefficients do model the dynamic nature of a signal but not necessarily the generalized context-independent signal likelihoods will give us more information about the dependence of the stream weights on the different types

of contexts.



REFERENCES

- [Boc93] E.L. Bocchieri and J.G. Wilson, "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-88-112
- Moore, R., "Evaluation of LPC Spectral Measures for Phonetic Unit Recognition", Report, Computer Science Department, Carnegie Mellon University, 1986
- Shibata, T., "An LQ Adaptive Speech Recognition System", Report, CMU-88-112

4. FUTURE WORK

So far we have only performed experiments with continuous 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, then there is no guarantee that the function will be convex. The highest probability path is not necessarily the one that maximizes the score for the parameters to increase the probability of the path. (The differentiation will be used to decrease the probability of the path.) (The differentiation will be used to decrease the probability of the path.) (The differentiation will be used to decrease the probability of the path.)

correct state C at time t , and let $LP_t(\alpha, C) := -\log P(x_t|C) = \sum_{i=1}^n c_i(t) \cdot \alpha_i(C)$ and $LP_t(\alpha, B) = \sum_{i=1}^n c_i(t) \cdot \alpha_i(B)$. Let $b_i(t)$ be the contribution of the stream i to the score for the best state B at time t , $b_i(t) = \frac{\partial LP_t(\alpha, B)}{\partial \alpha_i(B)}$. This means that $\sum_{i=1}^n c_i(t) \cdot \alpha_i(C) \leq \sum_{i=1}^n b_i(t) \cdot \alpha_i(B)$. (5)

of the training procedure is to find $\alpha_j(B)$ and $\alpha_j(C)$ such that $LP_t(\alpha, C)$ decreases and $LP_t(\alpha, B)$ increases. Here, we have ignored that the actual size of the infinitesimal step is somewhat greater than ϵ , resulting in a somewhat greater denominator. But $LP_t(\alpha, S)$ with respect to $\alpha_j(S)$. The update rule will then be $\alpha_j(S) \leftarrow \alpha_j(S) + \lambda \frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$.

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

$\frac{\partial LP_t(\alpha, B)}{\partial \alpha_j(B)} = b_j$ (4)

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 $\frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$ will not be independent of each other. Because of the above mentioned summation constraint, many gradient problems do not have a simple solution. For example, if we have a constraint $\sum_{i=1}^n \alpha_i(S) = 1$, then increasing $\alpha_j(S)$ by δ_j means that we must decrease some other $\alpha_i(S)$ by δ_j . This step definition needs the summation constraint and to maintain it unchanged. Other step

step from $\alpha_j(S) = (\alpha_1(S), \dots, \alpha_n(S))$ to $\alpha_j(S) = (\alpha_1(S) + \delta_1, \dots, \alpha_n(S) + \delta_n)$ could thus be defined as $\delta_j := \epsilon$ and $\delta_i := -\epsilon$ for $i \neq j$. This step definition

needs the summation constraint and to maintain it unchanged. Other step

