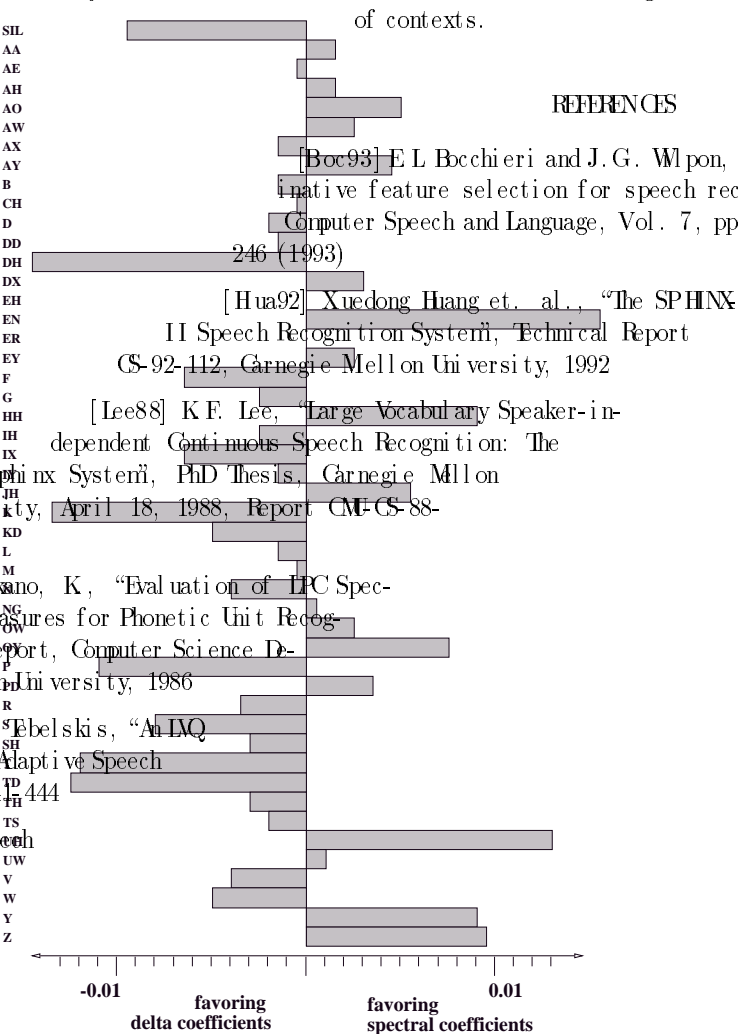


phoneme tend more to favoring delta-coefficients of features, like e.g. delta-spectral-coefficients, although one might expect that these coefficients' delta-delta-spectral-coefficients, power, acoustics are rather static and less context-dependent. Certainly, one fact that delta coefficients do model the dynamic delta-feature. Experiments with non-contextualized triphones including cross-word triphones 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] Xuedong Huang et. al., "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 CMU-CS-88-144

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

Joe Ebeliski, "An Adaptive Speech Recognition System", Report CMU-CS-88-144

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, then there is no guarantee that the function will be convex. The highest probability path is not necessarily the one that minimizes the loss function. The parameters to be updated are not necessarily the same as the parameters of the model. The step definition is not necessarily the same as the step definition of the model. The step definition is not necessarily the same as the step definition of the model.

correct state C at time t , and let $LP_t(\alpha, C) := -\log P(x_i|C)$, $\sum_{i=1}^n c_i(t) \cdot \alpha_i(C)$ and $LP_t(\alpha, B) = -\log P(x_i|B)$, $\sum_{i=1}^n b_i(t) \cdot \alpha_i(B)$. Let $b_i(t)$ be the contribution of the stream i to the score for the best state B at time t , $\sum_{i=1}^n b_i(t) \cdot \alpha_i(B)$. 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 $\alpha_j(S) \leftarrow \alpha_j(S) + \lambda \cdot \frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$.

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)

$\frac{d^* LP_t(\alpha, B)}{d \alpha_j(B)} = b_1 + \frac{\alpha_2(B) b_2}{d \alpha_j(B)}$ (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 $\frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$ will be positive for every update. This is not surprising because the gradient of the loss function is always positive. The partial derivative $\frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$ will be positive for every update. This is not surprising because the gradient of the loss function is always positive. The partial derivative $\frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$ will be positive for every update. This is not surprising because the gradient of the loss function is always positive.

step $\alpha_j(B) \leftarrow \alpha_j(B) + \lambda \cdot \frac{\partial LP_t(\alpha, B)}{\partial \alpha_j(B)}$ (feature F_j in S) $\alpha_j(S) \leftarrow \alpha_j(S) + \lambda \cdot \frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$ (feature F_j in S) $\alpha_j(S) \leftarrow \alpha_j(S) + \lambda \cdot \frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$ (feature F_j in S) $\alpha_j(S) \leftarrow \alpha_j(S) + \lambda \cdot \frac{\partial LP_t(\alpha, S)}{\partial \alpha_j(S)}$ (feature F_j in S)

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, PA 15213-3890, USA

1. TRAINING

ABSTRACT

which uses n information streams x_t , and for a given HMM state i , 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 the i -th stream x_t given the state i and the state sequence 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 the i -th stream x_t given the state i and the state sequence 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 the i -th stream x_t given the state i and the state sequence S .

(2)