

*Theorem:* Suppose  $\mu_1 > \mu_2$  (i.e., priority policy  $A$  is optimal if  $\mu_1$  and  $\mu_2$  are known). Then

$$\lim_{\substack{\epsilon, \delta \rightarrow 0 \\ N \rightarrow \infty}} R_A(\epsilon, \delta, N) = 1.$$

In other words, for any given  $0 < \beta < 1$ , there exist  $\epsilon_1, \delta_1$ , and  $N_1$  such that for all  $\epsilon < \epsilon_1, \delta < \delta_1$ , and  $N > N_1$ , it is assured that  $R_A(\epsilon, \delta, N) \geq 1 - \beta$ .

*Proof:* From (4)

$$\begin{aligned} \rho &= a_2/a_1 = \frac{\mu_1 \delta + o(\delta)}{\mu_2 \delta + o(\delta)} \\ &= \frac{\mu_1 + o(\delta)/\delta}{\mu_2 + o(\delta)/\delta} = \frac{\mu_1}{\mu_2} + O(\delta) \end{aligned}$$

and  $O(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ . This relation is easily verified by cross-multiplying. Also,

$$\rho^N = (\mu_2/\mu_1)^N + O(\delta).$$

Since  $\rho_1 = O(\epsilon)$ , from (5.2) and (5.4) for  $i = 1 \cdots N - 1, \pi_i = O(\epsilon)$  and

$$\pi_0^{-1} = 1 + (\mu_2/\mu_1)^N + O(\delta) + O(\epsilon).$$

Hence from (6)

$$\begin{aligned} R_A(\epsilon, \delta, N) &= \pi_0 + O(\epsilon) \\ &= \left[1 + (\mu_2/\mu_1)^N + O(\delta) + O(\epsilon)\right]^{-1} + O(\epsilon). \end{aligned}$$

Thus

$$\begin{aligned} \lim_{\substack{\epsilon, \delta \rightarrow 0 \\ N \rightarrow \infty}} R_A(\epsilon, \delta, N) &= \lim_{N \rightarrow \infty} \left[1 + (\mu_2/\mu_1)^N\right]^{-1} \\ &= 1, \quad \text{if } \mu_1 > \mu_2. \quad \text{Q.E.D.} \end{aligned}$$

### III. REMARKS, SIMULATIONS, AND CONCLUSION

In this work we demonstrated how a simple finite memory fixed-structure learning scheme may be effectively used to learn the priority assignment at a single server service station with two job-streams. The service times of jobs in each stream is exponentially distributed, whose parameters (and hence the optimal priority assignment policy) are unknown in the beginning. The algorithm parameters, in Section II, may be suitably chosen so that in steady state the probability of selecting the optimal action is arbitrarily close to one.

Simulation comparison of the fixed structure scheme and the variable structure scheme of [2] is presented in Table I. In the variable structure scheme of [2], the probability of selecting a policy is updated on the basis of the sample mean estimate gathered for each stream. This updating depends on a step size  $a$ . By selecting small values of  $a$  ( $a \rightarrow 0$ ), the steady-state probability of choosing the optimal policy may be made close to one as desired. The sample mean estimate for a stream is updated after serving each job from that stream. In Table I, each simulation point is averaged over 500 runs. The fixed-structure scheme is faster on the average, and it also has a low variance, during convergence. The reasons for low variance may be found in the observation that state transitions occur only when events that provide maximum evidence of one hypothesis against the other occur. In general, updating events [2] that are not maximum information events may initially lead the learning algorithm in a wrong direction. In such a case it takes on the average longer to converge to the optimum, resulting in slower convergence. For example, suppose  $\mu_1 > \mu_2$ , but the arrival rate ( $\lambda_2$ ), for stream 2 is much larger than ( $\lambda_1$ ) than for stream 1. Suppose a learning algorithm derives its decision on the basis of sample means (for each stream), which is updated after completion of every job.

TABLE I  
FIXED VERSUS VARIABLE STRUCTURE SCHEMES

$t$ (seconds)	Fixed $\epsilon = 0.05, \delta = 0.1, N = 20$		Variable Structure Scheme [2] Step Size $a = 0.1$	
	$P_A(t)$	$\sigma_A^2(t)$	$P_A(t)$	$\sigma_A^2(t)$
0 (start)	0.5	0	0.5	0
50	0.446	0.220	0.472	0.28
100	0.412	0.182	0.435	0.241
150	0.364	0.141	0.392	0.210
200	0.305	0.095	0.354	0.177
250	0.293	0.095	0.323	0.149
300	0.187	0.059	0.291	0.121
350	0.133	0.032	0.258	0.095
400	0.102	0.016	0.226	0.068
450	0.079	0.008	0.201	0.042
500	0.058	0.005	0.183	0.027
600	0.054	0.002	0.164	0.018
700	0.054	0.002	0.149	0.009
800	0.053	0.001	0.126	0.006
900	0.053	0.001	0.105	0.005
1500	0.052	0.0004	0.061	0.001

$t \geq 0$  time in seconds.

$\lambda_i$  Poisson arrival rate (jobs/seconds) for stream  $i = 1, 2$ .

$1/\mu_i$  average service time in seconds for stream  $i = 1, 2$ .

$\lambda_i = 0.25, \lambda_2 = 0.75, \mu_2 = 2\mu_1 = 2$ .

$P_A(t)$  average probability of choosing optimum policy  $A$  at time  $t$ .

$\sigma_A^2(t)$  variance of probability of choosing  $A$  at time  $t$ .

Since  $\lambda_2 > \lambda_1$ , the estimate of the mean converges near its true value  $\mu_2$  faster than that for stream 1. If, however, at this time the noisy estimate of  $\mu_1$  is smaller, then the algorithm would have drifted near the nonoptimal policy. Hence such a scheme would take a longer time to converge to the optimum.

The fixed-structure scheme proposed here bears resemblance to the optimal finite memory scheme of Cover-Hellman [4] for the two-armed bandit problem. Optimal finite memory schemes for the multiaction (greater than two) learning problems are not known in the literature. Extending the results of this paper to the case of more than two job streams, and relaxing assumption of the exponential service time to arbitrary distributions, are subjects of further investigation.

### REFERENCES

- [1] K. S. Narendra and M. A. L. Thatchachar, "Learning automata—A survey," *IEEE Trans. Syst., Man, Cybern.*, pp. 323-334, 1984.
- [2] M. R. Meybodi and S. Laxmivarahan, "A learning approach to priority assignment in a two class M/M/1 queueing system with unknown parameters," in *Proc. Yale Workshop on Adaptive System Theory*, June 1983, pp. 106-109.
- [3] A. Cobham, "Priority assignment in waiting line problems," *Operations Res.*, vol. 2, pp. 70-76, 1954.
- [4] T. Cover and M. Hellman, "The two-armed-bandit problem with time invariant finite memory," *IEEE Trans. Inform. Theory*, vol. IT-6, pp. 185-195, 1970.

### Fuzzy Set Theoretic Measure for Automatic Feature Evaluation

SANKAR K. PAL, SENIOR MEMBER, IEEE, AND BASABI CHAKRABORTY, STUDENT MEMBER, IEEE

*Abstract*—The terms *index of fuzziness*, *entropy*, and  *$\pi$ -ness*, which give measures of fuzziness in a set, are used to define an index of feature

Manuscript received April 5, 1984; revised August 27, 1985 and April 21, 1986.

The authors are with the Electronics and Communication Science Unit, Indian Statistical Institute, 203 Barrackpore Trunk Road, Calcutta 700035, India.

IEEE Log Number 8609835.

evaluation in pattern recognition problems in terms of their intraclass and interclass measures. The index value decreases as the reliability of a feature in characterizing and discriminating different classes increases. The algorithm developed has been implemented in cases of vowel and plosive identification problem using formant frequencies and different  $S$  and  $\pi$  membership functions.

## I. INTRODUCTION

The process of selecting the necessary information to present to the decision rule is called *feature selection*. Its main objective is to retain the optimum salient characteristics necessary for the recognition process and to reduce the dimensionality of the measurement space so that effective and easily computable algorithms can be devised for efficient classification.

The criterion of a good feature is that it should be unchanging with any other possible variation within a class, while emphasizing differences that are important in discriminating between patterns of different types. One of the useful techniques to achieve this is clustering transformation [1]–[3], which maximizes/minimizes the intersets/intraset distance using a diagonal transformation, such that smaller weights are given to features having larger variance (less reliable). Other separability measures based on information theoretic approach include divergence, Bhattacharyya coefficient, and the Kolmogorov variational distance [1]–[7].

The present work demonstrates an application of the theory of fuzzy sets to the problem of evaluating feature quality. The terms *index of fuzziness* [8], *entropy* [9], and  *$\pi$ -ness* [10] provide measures of fuzziness in a set and are used here to define the measure of separability in terms of their interclass and intraclass measurements. These two types of measurements are found to reflect the concept of intersets and intraset distances in classical set theory. An index of feature evaluation is then defined using these measures such that the lower the value of the index for a feature, the greater is the importance (quality) of the feature in recognizing and separating classes in the feature space.

It is also to be mentioned here that the above parameters provide algorithms for automatic segmentation [11] of grey tone image and measuring enhancement quality [12] of an image.

Effectiveness of the algorithm is demonstrated on vowel, plosive consonant, and speaker recognition problems using formant frequencies and their different combinations as feature set and  $S$  and  $\pi$  functions [13]–[15] as membership functions.

## II. FUZZY SETS AND MEASUREMENTS OF FUZZINESS

### A. Fuzzy Sets

A fuzzy set  $A$  with its finite number of supports  $x_1, x_2, \dots, x_n$  in the universe of discourse  $U$  is formally defined as

$$A = \{(\mu_A(x_i), x_i)\}, \quad i = 1, 2, \dots, n \quad (1)$$

where the characteristic function  $\mu_A(x_i)$  known as membership function and having positive value in the interval  $[0, 1]$  denotes the degree to which an event  $x_i$  may be a member of  $A$ . A point  $x_i$  for which  $\mu_A(x_i) = 0.5$  is said to be a crossover point of the fuzzy set  $A$ .

Let us now give some measures of fuzziness of a set  $A$ . These measures define, on a global sense, the degree of difficulty (ambiguity) in deciding whether an element  $x_i$  would be considered as a member of  $A$ .

### B. Index of Fuzziness

The index of fuzziness  $\gamma$  of a fuzzy set  $A$  having  $n$  supporting points reflects the degree of ambiguity present in it by measuring the distance between  $A$  and its nearest ordinary set  $\tilde{A}$  and is defined as [8]

$$\gamma(A) = \frac{2}{n^{1/k}} d(A, \tilde{A}) \quad (2)$$

where  $d(A, \tilde{A})$  denotes the distance between  $A$  and its nearest

ordinary set  $\tilde{A}$ . The set  $\tilde{A}$  is such that

$$\mu_{\tilde{A}}(x_i) = 0, \quad \text{if } \mu_A(x_i) \leq 0.5 \quad (3a)$$

and

$$\mu_{\tilde{A}}(x_i) = 1, \quad \text{if } \mu_A(x_i) > 0.5. \quad (3b)$$

The positive constant  $k$  appears in order to make  $\gamma(A)$  lie between zero and one, and its value depends on the type of distance function used. For example,  $k = 1$  for a generalized Hamming distance, whereas  $k = 2$  for a Euclidean distance. The corresponding indices of fuzziness are called the linear index of fuzziness  $\gamma_l(A)$  and the quadratic index of fuzziness  $\gamma_q(A)$ . Considering  $d$  to be a generalized Hamming distance, we have

$$\begin{aligned} d_l(A, \tilde{A}) &= \sum_i |\mu_A(x_i) - \mu_{\tilde{A}}(x_i)| \\ &= \sum_i \mu_{A \cap \bar{A}}(x_i), \quad i = 1, 2, \dots, n \end{aligned} \quad (4)$$

and

$$\gamma_l(A) = \frac{2}{n} \sum_i \mu_{A \cap \bar{A}}(x_i), \quad i = 1, 2, \dots, n \quad (5)$$

where  $\mu_{A \cap \bar{A}}(x_i)$  denotes the membership of  $x_i$  to a set which is the intersection of the fuzzy set  $A$  and its complement  $\bar{A}$  and is defined as

$$\mu_{A \cap \bar{A}}(x_i) = \min \{ \mu_A(x_i), (1 - \mu_A(x_i)) \}, \quad i = 1, 2, \dots, n.$$

Considering  $d$  to be an Euclidean distance, we have

$$\gamma_q(A) = \frac{2}{\sqrt{n}} \left[ \sum_i (\mu_A(x_i) - \mu_{\tilde{A}}(x_i))^2 \right]^{1/2}, \quad i = 1, 2, \dots, n. \quad (6)$$

### C. Entropy

The term entropy of a fuzzy set  $A$  is defined according to Deluca and Termini [9] as

$$H(A) = \frac{1}{n \ln 2} \sum_i S_n(\mu_A(x_i)), \quad i = 1, 2, \dots, n \quad (7)$$

with

$$\begin{aligned} S_n(\mu_A(x_i)) &= -\mu_A(x_i) \ln(\mu_A(x_i)) \\ &\quad - (1 - \mu_A(x_i)) \ln(1 - \mu_A(x_i)). \end{aligned} \quad (8)$$

In (7) and (8),  $\ln$  stands for natural logarithm (i.e., base  $e$ ). However, any other base would serve the purpose because of the normalization factor  $\ln 2$  in (7).

$\gamma(A)$  and  $H(A)$  are such that (from (5)–(7))

$$\gamma_{\min} = H_{\min} = 0(\text{min}), \quad \text{for } \mu_A(x_i) = 0 \text{ or } 1 \quad \text{for all } i \quad (9a)$$

$$\gamma_{\max} = H_{\max} = 1(\text{max}), \quad \text{for } \mu_A(x_i) = 0.5, \quad \text{for all } i. \quad (9b)$$

Suppose  $\mu_A(x_i) = 0.5$ , for all  $i$ . Then  $\mu_{\tilde{A}}(x_i) = 0$ , for all  $i$ , and

$$\gamma_l(A) = \frac{2}{n} \sum_i \left( \frac{1}{2} \right) = \frac{2}{n} \cdot \frac{n}{2} = 1$$

$$\gamma_q(A) = \frac{2}{\sqrt{n}} \left( \sum_i \left( \frac{1}{2} \right)^2 \right)^{1/2} = \frac{2}{\sqrt{n}} \cdot \frac{\sqrt{n}}{2} = 1$$

and

$$H(A) = \frac{1}{n \ln 2} \sum_i \left( -\ln \frac{1}{2} \right) = \frac{1}{n \ln 2} \cdot n \ln 2 = 1.$$

Therefore,  $\gamma$  and  $H$  increase monotonically in the interval  $[0, 0.5]$

and decrease monotonically in  $[0.5, 1]$  with a maximum of one at  $\mu = 0.5$ .

#### D. $\pi$ -ness

The  $\pi$ -ness of  $A$  is defined as [10]

$$\pi(A) = \frac{1}{n} \sum_i G_\pi(x_i), \quad i = 1, 2, \dots, n \quad (10)$$

where  $G_\pi$  is any  $\pi$  function as explained in the Section III.

$G_\pi$  ( $0 \leq G_\pi \leq 1$ ) increases monotonically in  $[x_i = 0$  to  $x_i = x_{\max}/2$ , say] and then decreases monotonically in  $[x_{\max}/2, x_{\max}]$  with a maximum of unity at  $x_{\max}/2$ , where  $x_{\max}$  denotes the maximum value of  $x_i$ .

### III. MEMBERSHIP FUNCTIONS

Let us now consider different  $S$  and  $\pi$  functions to obtain  $\mu_A(x_i)$  from  $x_i$ . The standard  $S$  function as defined by Zadeh [13] has the form

$$\mu_{AS}(x_i; a, b, c) = 0, \quad x_i \leq a \quad (11a)$$

$$= 2[(x_i - a)/(c - a)]^2, \quad a \leq x_i \leq b \quad (11b)$$

$$= 1 - 2[(x_i - c)/(c - a)]^2, \quad b \leq x_i \leq c \quad (11c)$$

$$= 1, \quad x_i \geq c \quad (11d)$$

in the interval  $[a, c]$  with  $b = (a + c)/2$ . The parameter  $b$  is known as the crossover point for which  $\mu_{AS}(b) = S(b; a, b, c) = 0.5$ .

Similarly, the standard  $\pi$  function has the form

$$\mu_{A\pi}(x_i; a, c, a') = \mu_{AS}(x_i; a, b, c), \quad x_i \leq c \quad (12a)$$

$$= 1 - \mu_{AS}(x_i; c, b', a'), \quad x_i \geq c \quad (12b)$$

in the interval  $[a, a']$  with  $c = (a + a')/2$ ,  $b = (a + c)/2$ , and  $b' = (a' + c)/2$ .  $b$  and  $b'$  are the crossover points, i.e.,  $\mu_{A\pi}(b) = \mu_{A\pi}(b') = 0.5$ , and  $c$  is the central point at which  $\mu_{A\pi} = 1$ .

Instead of using the standard  $S$  and  $\pi$  functions one can also consider the following equation as defined by Pal and Dutta Majumder [14], [15]

$$\mu_A(x_i) = G(x_i) = \left[ 1 + \left( \frac{|\hat{x}_n - x_i|}{F_d} \right) \right]^{-F_e} \quad (13)$$

which approximates the standard membership functions.

$F_e$  and  $F_d$  (two positive constants) are known respectively as exponential and denominational fuzzy generators and control the crossover point, bandwidth, and hence the symmetry of the curve about the crossover point.  $\hat{x}_n$  is the reference constant such that the function represents an  $S$ -type function  $G_S$  for  $\hat{x}_n = x_{\max}$  and a  $\pi$ -type function  $G_\pi$  for  $\hat{x}_n = x_i$ ,  $0 < x_i < x_{\max}$ , where  $x_{\max}$  represents the maximum value of  $x_i$ .

### IV. FEATURE EVALUATION INDEX

Let  $C_1, C_2, \dots, C_j, \dots, C_m$  be the  $m$ -pattern classes in an  $N$ -dimensional  $(X_1, X_2, \dots, X_q, \dots, X_N)$  feature space  $Q_X$ . Also, let  $n_j$  ( $j = 1, 2, \dots, m$ ) be the number of samples available from class  $C_j$ . The algorithms for computing  $\gamma$ ,  $H$ , and  $\pi$ -ness values of the classes in order to provide a quantitative index for feature evaluation are described in this section.

#### A. Computation of $\gamma$ and $H$ Using Standard $\pi$ Function

Let us consider the standard  $\pi$  function (12) for computing  $\gamma$  and  $H$  of  $C_j$  along the  $q$ th component and take the parameters of the function as

$$c = (x_{qj})_{av} \quad (14a)$$

$$b' = c + \max \left\{ |(x_{qj})_{av} - (x_{qj})_{\max}|, |(x_{qj})_{av} - (x_{qj})_{\min}| \right\} \quad (14b)$$

with

$$b = 2c - b' \quad (14c)$$

$$a = 2b - c \quad (14d)$$

$$a' = 2b' - c \quad (14e)$$

where  $(x_{qj})_{av}$ ,  $(x_{qj})_{\max}$ , and  $(x_{qj})_{\min}$  denote the mean, maximum, and minimum values respectively, computed along the  $q$ th coordinate axis, over all the  $n_j$  samples in  $C_j$ .

Since  $\mu(c) = \mu((x_{qj})_{av}) = 1$ , the values of  $\gamma$  and  $H$  are zero at  $c = (x_{qj})_{av}$  and would tend to unity (9) as we move away from  $c$  towards either  $b$  or  $b'$  of the  $\pi$  function (i.e., from mean towards boundary of  $C_j$ ). The lower the value of  $\gamma$  or  $H$  along the  $q$ th component in  $C_j$ , the greater would be the number of samples having  $\mu(x) \approx 1$  (or, the less would be the difficulty in deciding whether an element  $x$  can be considered, on the basis of its  $q$ th measurement, a member of  $C_j$  or not) and hence the greater would be the tendency of the samples to cluster around its mean value, resulting in less internal scatter or less intraset distance or more compactness of the samples along the  $q$ th axis within  $C_j$ . Therefore, the reliability (goodness) of a feature in characterizing a class increases as its corresponding  $\gamma$  or  $H$  value within the class (computed with  $\pi$  function) decreases.

The value of  $\gamma$  or  $H$  thus obtained along the  $q$ th coordinate axis in  $C_j$  may be denoted by  $\gamma_{qj}^\pi$  or  $H_{qj}^\pi$ .

Let us now pool together the classes  $C_j$  and  $C_k$  ( $j, k = 1, 2, \dots, m, j \neq k$ ) and compute the mean  $(x_{qjk})_{av}$ , maximum  $(x_{qjk})_{\max}$  and minimum  $(x_{qjk})_{\min}$  values of the  $q$ th component over all the samples (numbering  $n_j + n_k$ ). The value of  $\gamma$  or  $H$  so computed with (14) would therefore increase as the goodness of the  $q$ th feature in discriminating pattern classes  $C_j$  and  $C_k$  increases, because there would be fewer samples around the mean  $(x_{qjk})_{av}$  of the combined class, resulting in  $\gamma$  or  $H \approx 0$ , and more samples far from the  $(x_{qjk})_{av}$ , giving  $\gamma$  or  $H \approx 1$ . Let us denote the  $\gamma$  and  $H$  value so computed by  $\gamma_{qjk}^\pi$  and  $H_{qjk}^\pi$ , which increase as the separation between  $C_j$  and  $C_k$  (i.e., separation between  $b$  and  $b'$ ) along the  $q$ th dimension increases or, in other words, as the steepness of  $\pi$  function decreases.

It is to be mentioned here that one can also replace  $(x_{qj})_{av}$ ,  $(x_{qj})_{\max}$  and  $(x_{qj})_{\min}$  of (14) by  $(x_{qjk})_{av}$ ,  $(x_{qj})_{av}$ , and  $(x_{qk})_{av}$ , respectively, to compute  $\gamma_{qjk}$  or  $H_{qjk}$ . In this case, only their absolute values but not their behavior, as described previously, would be affected.

#### B. Computation of $\gamma$ and $H$ Using Standard $S$ Function

For computing  $\gamma$  and  $H$  of  $C_j$  along the  $q$ th component let us now take the parameters of  $S$  function (11) as

$$b = (x_{qj})_{av} \quad (15a)$$

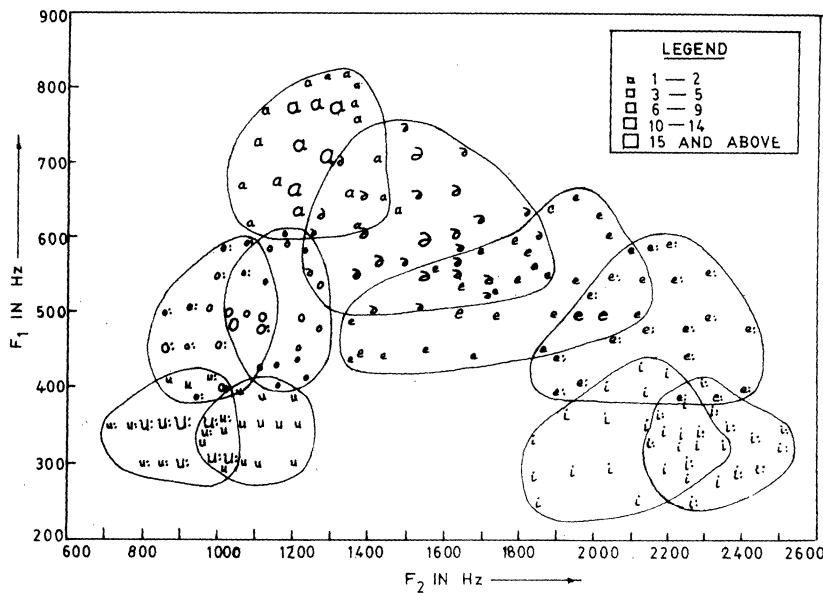
$$c = b + \max \left\{ |(x_{qj})_{av} - (x_{qj})_{\max}|, |(x_{qj})_{av} - (x_{qj})_{\min}| \right\} \quad (15b)$$

where

$$a = 2b - c. \quad (15c)$$

Since  $\mu(b) = \mu((x_{qj})_{av}) = 0.5$ , the values of  $\gamma$  and  $H$  are 1 at  $b = (x_{qj})_{av}$  and would tend to zero (9) as we move away from  $b$  towards either  $c$  or  $a$  of the  $S$  function. The higher the value of  $\gamma$  or  $H$ , the greater would be the number of samples having  $\mu(x) \approx 0.5$  and hence the greater would be the tendency of the samples to cluster around its mean value, resulting in less internal scatter within the class. Therefore, unlike the case with  $\pi$  function, the reliability (goodness) of a feature in characterizing a class  $C_j$  increases as its corresponding  $\gamma_{qj}^s$  or  $H_{qj}^s$  value within the class increases.

Similarly, if we now pool together the classes  $C_j$  and  $C_k$  ( $j, k = 1, 2, \dots, m, j \neq k$ ) and compute the mean, maximum and minimum values of the  $q$ th component over all the samples ( $n_j + n_k$ ), then the value of  $\gamma_{qjk}^s$  or  $H_{qjk}^s$  so computed with (15) would therefore decrease as the goodness of the  $q$ th feature in

Fig. 1. Vowel diagram in  $F_1$ - $F_2$  plane.

discriminating pattern classes  $C_j$  and  $C_k$  increases; because there would be fewer samples around the mean of the classes  $C_j$  and  $C_k$ , resulting in  $\gamma$  or  $H \approx 1$ , and more samples far from the mean, giving  $\gamma$ , or  $H \approx 0$ .

It therefore appears that the  $\gamma_{qj}$  (or  $H_{qj}$ ) and  $\gamma_{qjk}$  (or  $H_{qjk}$ ) reflect the concept of intraset and interset distances, respectively, in a classical feature-selection problem. With decrease in intraset and interset distances along  $q$ th component in  $C_j$ , the values of  $\gamma_{qj}$  (or  $H_{qj}$ ) and  $\gamma_{qjk}$  (or  $H_{qjk}$ ) are seen to decrease, or increase, when computed using the  $\pi$ , or  $S$  function.

### C. Computation of $\pi$ -ness

Similarly, for computing  $\pi_{qj}$  along the  $q$ th dimension in  $C_j$ , the parameters of the  $\pi$  function are set as follows:

$$c = (x_{qj})_{av} \quad (16a)$$

$$a' = c + \max \left\{ |(x_{qj})_{av} - (x_{qj})_{\max}|, |(x_{qj})_{av} - (x_{qj})_{\min}| \right\} \quad (16b)$$

with

$$a = 2c - a', \quad b = (a + c)/2, \quad b' = (a' + c)/2. \quad (16c)$$

For computing  $\pi_{qjk}$ , the classes  $C_j$  and  $C_k$  are pooled together and these parameters are obtained from  $(n_j + n_k)$  samples. Like the  $\gamma^s$  (or  $H^s$ ) value obtained with  $S$  function,  $\pi_{qj}$  and  $\pi_{qjk}$  increase as intraset and interset distances in  $C_j$  decrease.

Considering these intra-class and inter-class measures in each case, the problem of evaluating feature quality in  $Q_X$  therefore reduces to minimizing/maximizing the values of

$$\gamma_{qj}^{\pi} \text{ or } H_{qj}^{\pi} / \gamma_{qj}^s \text{ or } H_{qj}^s \text{ or } \pi_{qj}$$

while maximizing/minimizing the values of

$$\gamma_{qjk}^{\pi} \text{ or } H_{qjk}^{\pi} / \gamma_{qjk}^s \text{ or } H_{qjk}^s \text{ or } \pi_{qjk}.$$

The feature-evaluation index for the  $q$ th feature is accordingly defined as

$$(FEI)_q = \frac{d_{qj} + d_{qk}}{d_{qjk}}, \quad j, k = 1, 2, \dots, m, \quad j \neq k, \quad q = 1, 2, \dots, N \quad (17a)$$

where  $d$  stands for  $\gamma^{\pi}$  or  $H^{\pi}$  and

$$(FEI)_q = \frac{d_{qjk}}{d_{qj} + d_{qk}} \quad (17b)$$

where  $d$  stands for  $\gamma^s$  or  $H^s$  or  $\pi$ -ness. The lower the value of  $(FEI)_q$ , the higher is, therefore, the quality (importance) of the  $q$ th feature in characterizing and discriminating different classes in  $Q_X$ .

## V. IMPLEMENTATION AND RESULTS

For implementation of the above algorithm, the test material is prepared from a set of nearly 600 discrete phonetically balanced speech units in consonant-vowel-consonant Telugu (a major Indian Language) vocabulary uttered by three male speakers in the age group of 30-35 years.

For vowel sounds of ten classes ( $\delta$ ,  $a$ ,  $i$ ,  $i:$ ,  $u$ ,  $u:$ ,  $e$ ,  $e:$ ,  $o$  and  $o:$ ) including shorter and longer categories, the first three formant frequencies at the steady state ( $F_1$ ,  $F_2$ , and  $F_3$ ) are obtained through spectrum analysis.

For consonants, eight unaspirated plosive sounds namely the velars  $/k, g/$ , the alveolars  $/t, d/$ , the dentals  $/t, d/$ , and the bilabials  $/p, b/$  in combination with six vowel groups ( $\delta$ ,  $a$ ,  $E$ ,  $I$ ,  $O$ ,  $U$ ) are selected. The formant frequencies are measured at the initial and the final state of the plosives. The details of processing and formant extraction are available in [14]-[16].

### A. Vowel Recognition

A set of 496 vowel sounds of ten different classes are used here as the data set with  $F_1$ ,  $F_2$ , and  $F_3$  as the features. Fig. 1 shows the feature space of vowels corresponding to  $F_1$  and  $F_2$  when longer and shorter categories are treated separately.

Fig. 2 shows the order of importance of formants in recognizing and discriminating different vowels as obtained with intra-class measures (diagonal cells) and FEI values (off-diagonal cells). Results using only  $S$  function in computing  $\gamma_i$  and  $H$  values are shown here. Lower triangular part of the matrix corresponds to the results obtained with standard  $S$  and  $\pi$  functions ((11) and (12)) whereas, the upper triangular portion gives the results corresponding to their approximated versions (13). While using (13) we selected the parameters as follows.

For  $S$ -type Function:

$$\hat{x}_n = (x_{qj})_{\max}, \quad \text{for computing } \gamma_{qj} \text{ or } H_{qj} \quad (18a)$$

$$= (x_{qjk})_{\max}, \quad \text{for computing } \gamma_{qjk} \text{ or } H_{qjk}. \quad (18b)$$

$F_e$  and  $F_d$  were selected in such a way that  $(x_q)_{av}$  corresponds to the crossover point, i.e.,  $G_S((x_q)_{av}) = 0.5$ . To keep the crossover point fixed at  $(x_q)_{av}$ , different values of  $F_e$  and  $F_d$  may be used to result in various slopes of  $S$  function.

	$\delta$	Q:	I	U	E	O	
$\delta$	2 2 2 2	1 1 2	2 2	1 5 2	2 4 2	2 2	
Q:	4 2 2	2 2 2	2 2 1	1 1 1	2 2 2	2 2 1	
I	2 2 2	2 2 2	1 1 2	2 2 4	1 1 2	2 2 2	
U	1 1 1	1 1 1	2 2 2	2 2 2	2 1 2	1 1 3	
E	2 2 2	2 2 2	1 1 1	2 2 2	3 5 5	1 1 1	4 4 2
O	2 2 2	2 2 2	2 2 2	1 1 1	4 4 4	2 2 2	4 2 1

Code	Importance of Features
1	$F_1 > F_2 > F_3$
2	$F_2 > F_1 > F_3$
3	$F_1 > F_3 > F_2$
4	$F_2 > F_3 > F_1$
5	$F_3 > F_1 > F_2$
6	$F_3 > F_2 > F_1$

Diagonal cell: Intra class Ambiguity.  
Off-diagonal cell: FEI.  
Left triangle: Standard membership function.  
Right triangle: Approximated version.  
Upper entry :  $\delta$   
Middle entry : H  
Lower entry :  $\pi$ -ness

Fig. 2. Order of importance of features.

For  $\pi$ -type Function:

$$\hat{x}_n = (x_{qj})_{av}, \quad \text{for computing } \pi_{qj} \quad (19a)$$

$$= (x_{qjk}), \quad \text{for computing } \pi_{qjk}. \quad (19b)$$

As crossover points have no importance here in measuring  $\pi$ -ness, the selection of fuzzifiers is not crucial.

The results using (13) (as shown in the lower triangular part of Fig. 2) were obtained for  $F_e = 1/16, 1/8, 1/4, 1/2$ , and 1 with the crossover point at  $(x_{qj})_{av}$ . These values of  $F_e$  were also found to yield optimum recognition score in earlier investigations on vowel and plosive identification [14], [15]. For computing  $\pi_{qj}$  and  $\pi_{qjk}$ ,  $F_d$  was selected to be 50 for  $F_e = 1/16$ .

Again, the order of importance as shown in Fig. 2 was obtained after pooling together the shorter and longer counterparts (differing mainly in duration) of a vowel. In a part of the experiment the shorter and longer categories were treated separately, and the order of importance of formants for the corresponding  $\gamma_{qj}$ ,  $H_{qj}$ , and  $\pi_{qj}$  values (intra-class measures) is listed in Table I. This is included for comparison with the diagonal entries of Fig. 2.

For vowel recognition (except for /E/, as shown from Fig. 2) the first two formants are found to be much more important than  $F_3$  (which is mainly responsible for speaker identification). Furthermore, better result has been obtained for the cases when the shorter and longer categories are pooled together than the cases when they are treated separately. The result agrees well with previous investigation [14]. From the FEI measures of different pair of classes (off-diagonal cells of Fig. 2),  $F_1$  is seen to be more important than  $F_2$  in discriminating the class combinations /U, O/, /I, E/, /a, U/, and / $\delta$ , U/, i.e., between /front and front/ or /back and back/ vowels. For the other combinations, i.e., discriminating between /front and back/ vowels,  $F_2$  is found to be the strongest feature. The above findings can readily be verified from Fig. 1.

Typical FEI values for  $F_1$ ,  $F_2$ , and  $F_3$  are shown in Table II to illustrate the relative difference in importance among the formants in characterizing a class.

Similar investigations have also been made in case of speaker identification problem using the same data set (Fig. 1) and  $\{F_1, F_2, F_3, F_3 - F_2, F_3 - F_1, F_3/F_2, F_3/F_1\}$  as the feature set. FEI values have been computed for each of the three speakers

TABLE I  
INTRACLASS AMBIGUITIES FOR SHORT AND LONG VOWEL CLASSES

Membership Function	Vowel Class							
	i	i:	u	u:	e	e:	o	o:
Standard	2	1	2	3	4	1	2	2
	2	1	2	6	6	1	2	2
	1	1	2	6	6	1	2	2
Approximated version	3	1	2	1	5	1	2	1
	3	1	2	1	5	1	2	1
	1	1	2	1	1	1	2	1

1:  $F_1 F_2 F_3$  2:  $F_2 F_1 F_3$   
3:  $F_1 F_3 F_2$  4:  $F_2 F_3 F_1$   
5:  $F_3 F_1 F_2$  6:  $F_3 F_2 F_1$

individually for all the vowel classes. Contrary to the vowel recognition problem,  $F_3$  and its combinations were found here to yield lower FEI values, i.e., more important than  $F_1$  and  $F_2$ —resembling well the earlier report [14].

### B. Plosive Recognition

A set of 588 unaspirated plosive consonants are used as the data set with  $\Delta F_1, \Delta F_2$  (the difference of the initial and final values of the first and second formants),  $\Delta T$  (duration),  $\Delta F_1/\Delta T, \Delta F_2/\Delta T$  (the rates of transition) as the feature set.

The order of importance of the features for plosive recognition according to FEI values does not seem to be very regular as has been obtained in case of vowel recognition problem. Here all five features have more or less importance in determining the plosive classes, contrary to the case of vowel recognition, where  $F_3$  has much less importance than  $F_1$  and  $F_2$  in defining the vowel classes. However, a qualitative assessment has been adopted here to formulate an idea about the quality of the features based on the measure of FEI.

Table III shows the number of times each feature has occupied a particular position of importance on the basis of FEI measure using  $\gamma, H$ , and  $\pi$ -ness values and different target vowels. Results corresponding to both standard membership functions and their approximated versions are included for comparison.

TABLE II  
TYPICAL FEI VALUES OF THE FORMANTS USING STANDARD MEMBERSHIP FUNCTIONS

Vowel Classes	FEI Values According to the Parameter								
	Index of Fuzziness			Entropy			$\pi$ -ness		
	$F_1$	$F_2$	$F_3$	$F_1$	$F_2$	$F_3$	$F_1$	$F_2$	$F_3$
$\delta, I$	0.4521	0.3378	0.7142	0.4829	0.3874	0.6226	0.4828	0.3778	0.6402
$a, U$	0.2093	0.3649	0.5970	0.3038	0.4300	0.5860	0.2755	0.4253	0.5963
$I, E$	0.4503	0.4591	0.5894	0.4822	0.4841	0.5581	0.4795	0.4846	0.5649
$I, O$	0.3133	0.1018	0.6329	0.3994	0.1811	0.5586	0.3826	0.1483	0.5645
$U, E$	0.4538	0.2696	0.5635	0.4839	0.3289	0.5282	0.4822	0.3158	0.5286
$U, O$	0.3126	0.5196	0.5720	0.3948	0.4999	0.5547	0.3859	0.4994	0.5630

TABLE III  
IMPORTANCE OF PLOSIVE FEATURES ACCORDING TO FEI

Membership Function	Position of Importance		
	First	Second	Third
Standard	$I_{36}, III_{34}, II_{20}$	$V_{42}, III_{25}, I_{17}$	$IV_{36}, I_{22}, II_{21}$
Approximate	$IV_{61}, V_{28}, I_{13}$	$III_{46}, I_{24}, V_{20}$	$III_{46}, I_{25}, II_{22}$
Standard	$I_{43}, III_{30}, II_{22}$	$III_{31}, II_{26}, I_{24}$	$II_{39}, III_{30}, V_{26}$
Approximate	$IV_{62}, I_{17}, II_{12}$	$III_{41}, I_{25}, IV_{17}$	$III_{44}, I_{24}, V_{21}$

1) Suffixes indicate the number of times the feature has occurred in the position.  
 2) I, II, III, IV, V represent the features  $\Delta F_1, \Delta F_2, \Delta T, \Delta F_1/\Delta T, \Delta F_2/\Delta T$ , respectively

Let us now consider the case of unvoiced plosive sounds with the standard membership function. The features  $\Delta F_1, \Delta T$ , and  $\Delta F_2$  were first in order of importance 36, 34, and 20 times, respectively. They occupied the second position 42, 25, and 17 times, respectively, and the third 36, 22, and 21 times, respectively. Considering the first three number of occurrences in first two positions and the first two number of occurrences in first two positions, it is seen that the set  $(\Delta F_1, \Delta T)$  is more effective than  $(\Delta F_2/\Delta T, \Delta F_2)$ , which is again more important than  $\Delta F_1/\Delta T$  in discriminating unvoiced plosive sounds. Similarly, the features  $(\Delta F_1, \Delta T, \text{ and } \Delta F_2)$  (particularly,  $\Delta F_1, \Delta T$ ) are seen to be more reliable than the others in characterizing voiced plosives using the standard membership function.

Let us now consider the cases of using approximate version of the membership functions (13). To discriminate unvoiced plosives the set  $(\Delta F_1/\Delta T, \Delta F_2/\Delta T)$  gives better characterizing feature than  $\Delta F_1$  and  $\Delta T$ ; whereas for the voiced counterparts  $(\Delta F_1/\Delta T, \Delta F_1, \Delta T)$  came out to be the best feature set.

From these discussions, the features  $\Delta F_1$  and  $\Delta T$  are overall found to be the most important in characterizing and discriminating different plosive sounds. The result conforms to the earlier findings [15], [16] obtained from the point of automatic plosive sound recognition.

VI. DISCUSSION AND CONCLUSION

An algorithm for automatic evaluation of feature quality in pattern recognition has been described using the terms *index of fuzziness*, *entropy*, and  $\pi$ -ness of a fuzzy set. These terms are used to define measures of separability between classes and compactness within a class when they are implemented with  $S$  and  $\pi$  membership functions. For example, when these measures are implemented with  $\pi$  function,  $d_{qj}^\pi$  ( $d$  stands for  $\gamma$  or  $H$ ) then decreases as compactness within  $j$ th class along  $q$ th direction increases and  $d_{qjk}^\pi$  increases as separability between  $C_j$  and  $C_k$  increases in the  $q$ th direction. If the classes  $C_j$  and  $C_k$  do not differ in mean value but differ in second order moment, i.e., variances are different, then

$$[(x_{qj})_{\max} - (x_{qj})_{\min}] > [(x_{qk})_{\max} - (x_{qk})_{\min}]$$

(assuming  $C_j$  with larger variance than  $C_k$ ). From (12) and (14)

we have

$$d_{qj}^\pi > d_{qk}^\pi,$$

i.e., the  $q$ th feature is more important in recognizing the  $k$ th class than the  $j$ th class. Value of the interset ambiguity  $d_{qjk}^\pi$  as expected, then decreases showing the deterioration in reliability (goodness) of the  $q$ th feature in discriminating  $C_j$  from  $C_k$ . Similar behavior would also be reflected for the third (representing skewness of a class) and higher order moments when they are different for  $C_j$  and  $C_k$  with the same mean value. The algorithm is found to provide satisfactory order of importance of the features in characterizing speech sounds, in discriminating different classes of speeches and also in identifying a speaker.

Since  $F_3$  and its higher formants ( $F_4, F_5, \dots$ ) are mostly responsible for identifying a speaker, we have considered in our experiment only  $F_3$  in addition to  $F_1$  and  $F_2$  for evaluating feature quality in vowel recognition problem.

It is to be mentioned here that the well-known statistical measures of feature evaluation such as Bhattacharyya coefficient, divergence, Kolmogorov variational distance, etc., theoretically take into account the interdependence of feature variables. Their computation involves multivariate numerical integration and estimation of probability density functions [4]. In practice in their computation, the features are usually treated individually to avoid computational difficulty [17]. Our proposed measure also treats the features individually. In fact, the present algorithm attempts to rank individual features according to their importance in characterizing and discriminating classes. Combination of features in doing so is not of interest. Furthermore, even in the case of independent feature, the algorithm is computationally more efficient than the aforesaid statistical measures.

ACKNOWLEDGMENT

The authors gratefully acknowledge Prof. D. Dutta Majumder for his interest in this work and Mrs. S. De Bhowmik and Mr. S. Chakraborty for preparing the manuscript.

REFERENCES

- [1] K. Fukunaga, *Introduction to Statistical Pattern Recognition and Machine Learning*. New York: Academic, 1972.
- [2] J. T. Tou and R. C. Gonzalez, *Pattern Recognition Principles*. Reading, MA: Addison-Wesley, 1974.
- [3] G. S. Sebestyen, *Decision Making Process in Pattern Recognition*. New York: Macmillan, 1962.
- [4] R. A. Devijver and J. Kittler, *Pattern Recognition—A Statistical Approach*. London: Prentice-Hall, 1982.
- [5] S. K. Pal and D. Dutta Majumder, *Fuzzy Mathematical Approach to Pattern Recognition*. New York: Wiley (Halsted Press), 1986.
- [6] A. Bhattacharyya, "On a measure of divergence between two multinomial populations," *Sankhya*, vol. 6, pp. 401-406, 1946.
- [7] T. Kailath, "The divergence and Bhattacharyya distance measures in signal detection," *IEEE Trans. Commun. Tech.*, vol. CT-15, pp. 52-60, 1967.
- [8] A. Kaufmann, *Introduction to the Theory of Fuzzy Subsets—Fundamental Theoretical Elements*, vol. 1. New York: Academic, 1975.
- [9] A. Deluca and S. Termini, "A definition of a nonprobabilistic entropy in the setting of fuzzy set theory," *Inform. Contr.*, vol. 20, pp. 301-302, 1972.
- [10] S. K. Pal, "Fuzzy set theoretic approach—A tool for speech and image recognition," in *Pattern Recognition Theory and Applications*. Amsterdam: D. Reidel, 1982, pp. 103-117.

- [11] S. K. Pal, R. A. King, and A. A. Hashim, "Automatic grey level thresholding through index of fuzziness and entropy," *Patt. Recognition Lett.*, vol. 1, pp. 141-146, Mar. 1983.
- [12] —, "A note on the quantitative measure of image enhancement through fuzziness," *IEEE Trans. Patt. Anal. Mach. Intell.*, vol. PAMI-4, pp. 204-208, 1982.
- [13] L. A. Zadeh, "Calculus of fuzzy restrictions," in *Fuzzy Sets and Their Application to Cognitive and Decision Process*. London: Academic, pp. 1-39, 1975.
- [14] S. K. Pal and D. Dutta Majumder, "Fuzzy sets and decision making approaches in vowel and speaker recognition," *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-7, pp. 625-629, 1977.
- [15] —, "On automatic plosive identification using fuzziness in property sets," *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-8, pp. 302-308, Apr. 1978.
- [16] A. K. Datta, N. R. Ganguli, and S. Ray, "Recognition of unaspirated plosives—A statistical approach," *IEEE Trans. Acoustics, Speech, Signal Processing*, vol. ASSP-28, pp. 85-91, Feb. 1980.
- [17] S. Ray, "The effectiveness of features in pattern recognition," Ph.D. thesis, Imperial College, Univ. of London, 1984.