

Serial Number                    09/285,173  
Filing Date                      18 March 1999  
Inventor                         Robert S. Lynch, Jr.  
                                       Peter K. Willett

NOTICE

The above identified patent application is available for licensing. Requests for information should be addressed to:

OFFICE OF NAVAL RESEARCH  
DEPARTMENT OF THE NAVY  
CODE 00CC  
ARLINGTON VA 22217-5660

**DISTRIBUTION STATEMENT A**  
Approved for Public Release  
Distribution Unlimited

DTIC QUALITY INSPECTED 4

19991007 047

1 Attorney Docket No. 79550

2

3 DATA REDUCTION SYSTEM FOR IMPROVING CLASSIFIER PERFORMANCE

4

5 STATEMENT OF GOVERNMENT INTEREST

6 The invention described herein may be manufactured and used  
7 by or for the Government of the United States of America for  
8 governmental purposes without the payment of any royalties  
9 thereon or therefor.

10

11 BACKGROUND OF THE INVENTION

12 (1) Field of the Invention

13 The invention relates to a data reduction system that  
14 reduces the dimensionality of neural network training data by  
15 finding features that most improve performance of the neural  
16 network.

17 (2) Description of the Prior Art

18 The use of classification systems to classify input data  
19 into one of several predetermined classes is well known. Their  
20 use has been adapted to a wide range applications including  
21 target identification, medical diagnosis, speech recognition,  
22 digital communications and quality control systems.

23 Classification of sonar signals into threats and non-threats  
24 is an important task for sonar operators. Neural networks have  
25 been proposed to help accomplish this task by receiving a signal

1 from the sonar system and analyzing characteristics of the signal  
2 for determining if the signal is originating from a vessel that  
3 is a military vessel that represents a threat or from a  
4 commercial vessel. Speed in making this determination is often  
5 of the essence.

6 Classification systems decide, given an input X, to which of  
7 several output classes X belongs. If known, measurable  
8 characteristics separate classes, the classification decision is  
9 straightforward. However, for most applications, such  
10 characteristics are unknown, and the classification system must  
11 decide which output class the input most closely resembles. In  
12 such applications, the output classes and their characteristics  
13 are modeled (estimated) using statistics for the classes derived  
14 from training data belonging to known classes. Thus, the  
15 standard classification approach is to first estimate the  
16 statistics from the given training data and then to apply a  
17 decision rule using these estimated statistics.

18 However, often there is insufficient training data to  
19 accurately infer the true statistics for the output classes which  
20 results in reduced classification performance or more occurrences  
21 of classification errors. Additionally, any new information that  
22 arrives with the input data is not combined with the training  
23 data to improve the estimates of the symbol probabilities.  
24 Furthermore, changes in symbol probabilities resulting from  
25 changes, which may be unobservable, in the source of test data,

1 the sensors gathering data or the environment often result in  
2 reduced classification performance. Therefore, if based on the  
3 training data, a classification system maintains a near zero  
4 probability for the occurrence of a symbol and the symbol begins  
5 to occur in the input data with increasing frequency,  
6 classification errors are likely to occur if the new data is not  
7 used in determining symbol probabilities.

8 Attempts to improve the classification performance and take  
9 advantage of information available in test data have explored  
10 combining the test data with the training data in modeling class  
11 statistics and making classification decisions. While these  
12 attempts have indicated that improved classification performance  
13 is possible, they have one or more drawbacks which limit or  
14 prevent their use for many classification systems.

15 The use of Bayesian classification is taught in the prior  
16 art for combining training data with test data is found in Merhav  
17 et al, "A Bayesian Classification Approach with Application to  
18 Speech Recognition," *IEEE Trans. Signal Processing*, vol. 39, no.  
19 10 (1991) pp. 2157-2166. In Merhav et al classification decision  
20 rules which depend on the available training and test data were  
21 explored. A first decision rule which is a Bayesian rule was  
22 identified. However, this classification rule was not fully  
23 developed or evaluated because the implementation and evaluation  
24 of the probability density functions required are extremely  
25 complex.



1 Another object of the invention is that such classification  
2 system should not include redundant and ineffectual data.

3 A further object of the invention is to provide a method for  
4 reducing feature vectors to only those values which affect the  
5 outcome of the classification.

6 Accordingly, this invention provides a data reduction method  
7 for a classification system using quantized feature vectors for  
8 each class with a plurality of features and levels. The  
9 reduction algorithm consisting of applying a Bayesian data  
10 reduction algorithm to the classification system for developing  
11 reduced feature vectors. Test data is then quantified into the  
12 reduced feature vectors. The reduced classification system is  
13 then tested using the quantized test data.

14 A Bayesian data reduction algorithm is further provided  
15 having by computing an initial probability of error for the  
16 classification system. Adjacent levels are merged for each  
17 feature in the quantized feature vectors. Level based  
18 probabilities of error are then calculated for these merged  
19 levels among the plurality of features. The system then selects  
20 and applies the merged adjacent levels having the minimum level  
21 based probability of error to create an intermediate  
22 classification system. Steps of merging, selecting and applying  
23 are performed until either the probability of error stops  
24 improving or the features and levels are incapable of further  
25 reduction.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

BRIEF DESCRIPTION OF THE DRAWINGS

A more complete understanding of the invention and many of the attendant advantages thereto will be readily appreciated as the same becomes better understood by reference to the following detailed description when considered in conjunction with the accompanying drawings wherein:

FIG. 1 is a functional block diagram showing implementation of the system of the current invention; and

FIG. 2 is a detailed diagram of the Bayesian data reduction algorithm of the current invention.

DESCRIPTION OF THE PREFERRED EMBODIMENT

The data reduction system is illustrated in the FIG. 1. This figure provides a functional block diagram of the major components of the data reduction system. Intended users of this system should implement it using FIG. 1, FIG. 2, and the associated formulas and algorithms described below, by writing a computer program in the language of their choice.

In step 10 of the data reduction system all training data for each class are represented as quantized feature vectors. The classification can have two or more classes. In the case when there are two classes, there are  $N_{target}$  quantized feature vectors for the *target* class and  $N_{nontarget}$  quantized feature vectors for the *nontarget* class. Each feature vector is quantized by mapping it to a symbol. There are  $M$  possible symbols representing the

1 number of discrete levels for a specific feature multiplied by  
2 the number of discrete levels for each feature. For example, a  
3 feature vector having three binary valued features can take on  
4 one of  $M = 8$  possible discrete symbols given by;  $(0, 0, 0)$ ,  $(0,$   
5  $0, 1), \dots, (1, 1, 1)$ .

6 In some cases, either one or all of the features will be  
7 continuous, and these features must then be discretized before  
8 the vectors are mapped to one of the  $M$  possible symbols. If a  
9 given set of thresholds does not exist for discretizing a  
10 particular feature then the feature should be discretized into a  
11 sufficient number of levels via percentiles. Ten discrete levels  
12 has been found to be adequate for most continuous features;  
13 however, other levels can be established depending on the  
14 sensitivity of the system to the feature vector and the capacity  
15 of the computer performing the data reduction. That is, to  
16 discretize a feature into ten levels its training data are used  
17 to define ten thresholds corresponding to ten percentile regions  
18 (e.g., the first threshold is found such that 10 percent of the  
19 data are less than it in value). This procedure is then repeated  
20 for the remaining continuous features. Notice also that there is  
21 no specified limit to the number of features used in the data  
22 reduction system. If the computational limits of the computer  
23 platform allow, using all known features is best. However, the  
24 same features must be used for each class, but it is not

1 necessary that the initial quantization of each feature be the  
2 same.

3 In step 12 the Bayesian data reduction algorithm is  
4 simultaneously applied to the quantized training data of all  
5 classes. The algorithm uses the Dirichlet distribution as a  
6 noninformative prior. The Dirichlet represents all symbol  
7 probabilities as uniformly-distributed over the positive unit-  
8 hyperplane. Using this prior, the algorithm works by reducing  
9 the quantization fineness,  $M$ , to a level which minimizes the  
10 average conditional probability of error,  $P(e)$ .

11 The formula for  $P(e)$  is the fundamental component of this  
12 algorithm, and it is given by

$$13 \quad P(e|X) = \sum_y \sum_x P(H_k) I_{\{z_k \leq z_l\}} f(y|x_k, H_k) + P(H_l) I_{\{z_k > z_l\}} f(y|x_l, H_l) \quad (1)$$

14 where, in the following  $k$  and  $l$  are exchangeable;

$$15 \quad z_k = f(y|x_k, H_k) = \frac{N_y!(N_k + M - 1)!}{(N_k + N_y + M - 1)!} \prod_{i=1}^M \frac{(x_{k,i} + y_i)!}{x_{k,i}! y_i!}; \quad (2)$$

16  $N$  is the number of feature vectors;

17  $k, l \in \{\text{target, nontarget}\}$ , and  $k \neq l$ ;

18  $M$  is the number of discrete symbols

19  $H_k$  is the hypothesis;

$$20 \quad H_k : p_y = p_k; \quad (3)$$

21  $p$  is the vector of probabilities;

$$22 \quad X \equiv (x_k, x_l); \quad (4)$$

1  $x_{k,i}$  is the number of symbol type  $i$  in the training data for class

2  $k$  and  $N_k \left\{ N_k = \sum_{i=1}^M x_{k,i} \right\};$

3  $y_i$  is the number of symbol type  $i$  in the test data and

4  $N_y \left\{ N_y = \sum_{i=1}^M y_i \right\};$

5  $f(\mathbf{y}|\mathbf{x}, H)$  is the probability distribution of  $\mathbf{y}$  which depends on the

6 parameter  $\mathbf{x}$  for the given hypothesis,  $H$ ; and

7  $I_{\{x\}}$  is the indicator function such that  $I_{\{x\}}=1$  when  $x$  is true and

8  $I_{\{x\}}=0$  when  $x$  is false.

9 For one test observation  $f(\mathbf{y}|\mathbf{x}_k, H_k)$  becomes

10 
$$f(\mathbf{y}|\mathbf{x}_k, H_k; N_y = 1) = \frac{x_{k,i} + 1}{N_k + M}, y_i = 1. \quad (5)$$

11 Given formula (1), the algorithm is implemented by using the  
12 following iterative steps as shown in FIG. 2.

13 In step 20, using the initial training data with  
14 quantization  $M$ , formula (1) is used to compute  $P(e|X; M)$ .  
15 In step 22, a feature is selected arbitrarily, and then a two  
16 adjacent levels of the feature are selected in step 24. Step 26  
17 merges the training data of those adjacent quantized symbols. In  
18 the binary case, quantized symbols containing a binary zero with  
19 are combined with those containing a binary one effectively  
20 removing the feature. In the continuous case, two levels are  
21 merged into one level removing the distinction between the two

1 levels. Step 28 uses the newly merged training data,  $X'$ , and the  
2 new quantization,  $M'$ , and again computes  $P(e|X'; M')$ . Step 30 is  
3 a loop wherein steps 22 through 28 are repeated for all adjacent  
4 feature quantizing levels, and all remaining features.

5 The algorithm then selects the merged configuration having  
6 the minimum probability of error,  $P(e|X'; M')$  in step 32 from the  
7 probabilities computed in step 28. The configuration with the  
8 minimum probability of error (or maximum probability of  
9 recognition) is then used as the new training data configuration  
10 for each class (i.e., the new quantization, and its associated  
11 discrete levels and thresholds for each feature). Step 34 is  
12 another loop which repeats steps 22 through 32 until the  
13 probability of error decreases no further, or until features can  
14 no longer be reduced, i.e.  $M' = 2$ .

15 In cases when the several probabilities are the same, the  
16 minimum can be selected arbitrarily. As an alternative the  
17 multiple configurations each having the same minimum  
18 probabilities can all be applied. By applying all  
19 configurations, computer processing time can be reduced at some  
20 increase in error. Accordingly, arbitrary selection of a single  
21 configuration is the preferred alternative.

22 Observe that the algorithm described above is "greedy" in  
23 that it chooses a best training data configuration at each  
24 iteration (see step 34 above) in the process of determining a  
25 best quantization fineness. A global search over all possible

1 merges and corresponding training data configurations may in some  
2 cases provide a lower probability of error at a higher  
3 computational cost. However, a simulation study involving  
4 hundreds of independent trials revealed that only about three  
5 percent of the time did the "greedy" approach shown above produce  
6 results different than a global approach. Additionally, the  
7 overall average probability of error for the two approaches  
8 differed by only an insignificant amount.

9       When the Bayesian data reduction algorithm finds the new  
10 quantization fineness upon completion of step 34 in FIG. 2, this  
11 new configuration can be established as in step 36. The  
12 resulting trained classifier can be tested as step 14 of FIG. 1.  
13 To test the classifier all test data from 16 are now quantized  
14 using the remaining features, and their associated discrete  
15 levels and threshold settings that were found in step 12 for the  
16 training data.

17       An advantage of the Bayesian data reduction algorithm of the  
18 current invention is that it permanently reduces, or eliminates,  
19 irrelevant and redundant features (as opposed to appropriately  
20 adjusting the weights of a neural network and keeping all  
21 features) from the training data. Thus, with the current  
22 invention features are important to correct classification are  
23 highlighted. With this, the algorithm presented here does not  
24 require the long training times that can accompany a neural  
25 network, nor does it require a randomized starting configuration.



1 Attorney Docket No. 79550

2

3 DATA REDUCTION SYSTEM FOR IMPROVING CLASSIFIER PERFORMANCE

4

5 ABSTRACT OF THE DISCLOSURE

6 A data reduction method for a classification system using  
7 quantized feature vectors for each class with a plurality of  
8 features and levels. The reduction algorithm consisting of  
9 applying a Bayesian data reduction algorithm to the  
10 classification system for developing reduced feature vectors.  
11 Test data is then quantified into the reduced feature vectors.  
12 The reduced classification system is then tested using the  
13 quantized test data.

14 A Bayesian data reduction algorithm is further provided  
15 having by computing an initial probability of error for the  
16 classification system. Adjacent levels are merged for each  
17 feature in the quantized feature vectors. Level based  
18 probabilities of error are then calculated for these merged  
19 levels among the plurality of features. The system then selects  
20 and applies the merged adjacent levels having the minimum level  
21 based probability of error to create an intermediate  
22 classification system. Steps of merging, selecting and applying  
23 are performed until either the probability of error stops  
24 improving or the features and levels are incapable of further  
25 reduction.

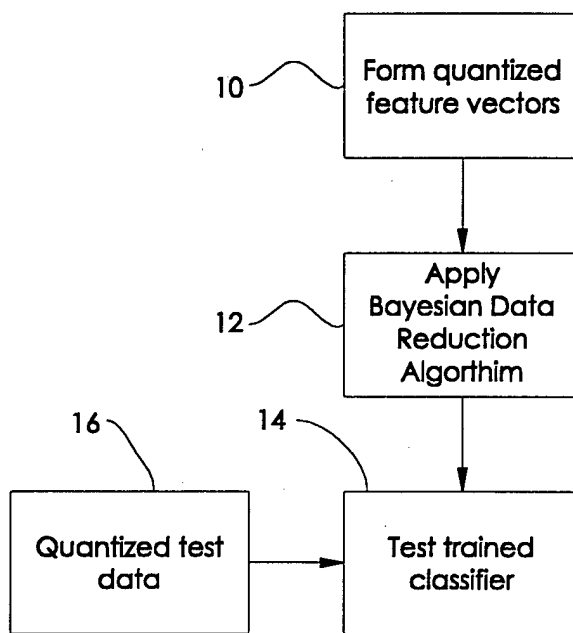


FIG. 1

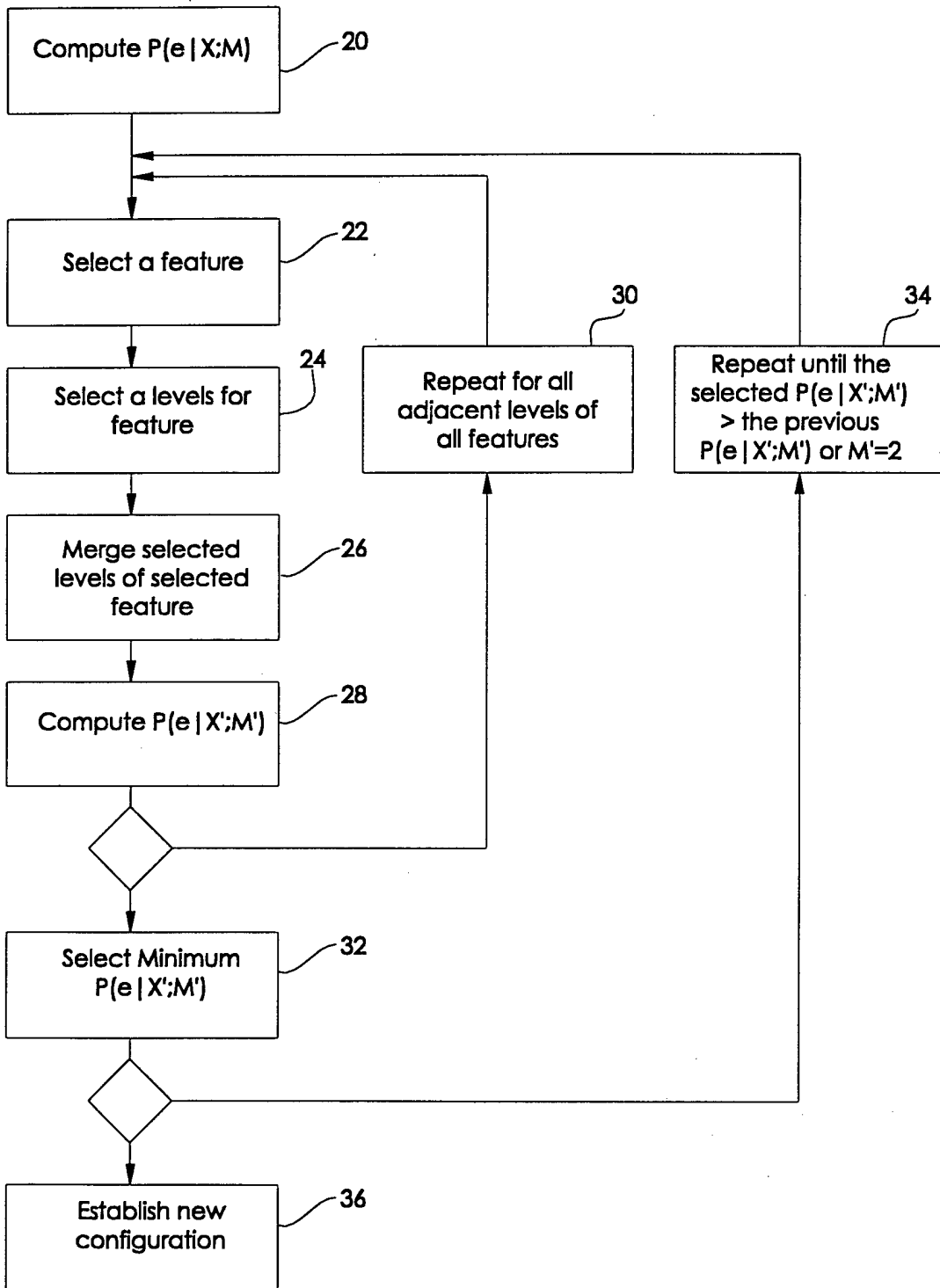


FIG. 2