

ORGANIZING MULTIPLE EXPERTS FOR EFFICIENT PATTERN RECOGNITION

By

Krassimir Ianakiev

July 2000

A DISSERTATION SUBMITTED TO THE
FACULTY OF THE GRADUATE SCHOOL OF STATE
UNIVERSITY OF NEW YORK AT BUFFALO
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

© Copyright 2000

by

Krassimir Ianakiev

Acknowledgments

I would like to thank my advisor Prof. V. Govindaraju for being a constant source of inspiration to me.

I would like to thank Prof. P. Scott, Prof S. Srihari and Dr. P. Slavik for taking the time to serve on my dissertation committee and for their suggested improvements to the dissertation.

Special thanks to Dr. A. Shekhawat.

Abstract

The purpose of this research is to investigate a general theory of classifier combination under complete abstraction of classifier details. In order to completely separate the design of classifiers and the design of combinator, we will treat each classifier as a black box whose sole interface to the combinator is a set of transformed decision vectors. Since the decision vector transformation is classifier dependent, it has to be provided. Also provided is a set of labeled data whose distribution reflects the environment in which the combined system will operate. All the necessary knowledge about classifiers can be derived through observations of their responses to the observation data. The classes to be recognized in this observation data set are also predetermined. This black box treatment of classifiers can be correlated, trained on different data sets, generating mixed types of outputs, and behaving in various ways. Under these presumptions, we wish to develop a combination theory.

Contents

Acknowledgments	iii
Abstract	iv
1 Introduction	1
1.1 Objective	3
1.2 Literature Review	5
1.3 Outline of the Dissertation	12
1.3.1 Measures	13
1.3.2 Applications	14
2 Measures	17
2.1 Deriving Pseudo-Probabilities of Correctness Given Scores (DPPS)	18

2.1.1	An Intuitive Problem Description	20
2.1.2	Theoretical Foundation of DPPS Methodology	22
2.1.3	DPPS Algorithm	25
2.1.4	Experimental Results	32
2.1.5	Some Final Remarks and Observations	33
2.2	Fuzzy Measure and Potential Improvement of Classifier Accuracy . .	34
2.2.1	Issue with traditional techniques	36
2.2.2	Our Approach	39
2.2.3	Topology	52
2.2.4	Fuzzy Sets	55
2.2.5	Bayesian View of the Membership Functions	57
2.3	Experiments	59
2.4	Summary	62
3	Applications-Classifier Combination	65
3.1	Confidence Combination Methods in Multi-Expert Systems	66
3.1.1	Accuracy Approximation	67

3.1.2	Optimal Combination Rule	69
3.1.3	Empirical Combination Rules	69
3.1.4	Optimal Rejection Threshold	73
3.2	Experimental Results	74
3.2.1	Conclusion	76
3.3	Fuzzy Control Structure in Multiple Parameter Systems	77
3.3.1	Handwritten Address Interpretation	78
3.3.2	Intuitive Understanding of Methodology	81
3.3.3	Topology, Measure and Membership of Fuzzy Sets	81
3.3.4	Experiments	85
3.4	Architecture for Classifier Combination Using Entropy Measure	86
3.4.1	Universal Combinator	89
3.4.2	Preliminary Results	99
3.4.3	Entropy Measure for Serial Classifier Combination	101
3.4.4	Greedy Algorithm	108
3.4.5	Experiments	110

3.4.6	Remarks	113
4	Conclusions	115

Chapter 1

Introduction

Pattern recognition applications require highly reliable classification which may be difficult for a single algorithm to achieve. Since there are a number of classification algorithms described in the literature, based on different theories and methodologies, a combination of these can be used to improve the classification performance. The idea is to take advantage of their strengths and avoid their weaknesses. While the idea is meritorious, it presents many challenges. Primarily, the decisions of the individual experts are often conflicting [9], [28], [31], [39], [38], [43], making the task non-trivial.

In general, combination of classification algorithms can be performed at two levels. First, combination is at the *feature* level. The features of each classifier are combined to form a joint feature vector in a new feature space and this newly formed vector is classified by methods applicable to the new feature space. Integration at this level avoids information loss. However, the potentially large dimensionality of the newly formed feature vector renders this approach impractical as lack of uniform representation of heterogeneous feature types can complicate the combination process.

A small set of features that can be discriminated by a simple method may require a complicated classification scheme when *other* sets of features are added. Thus selection of an appropriate classifier and its training process becomes difficult.

The second method of combination is at the classifier *decision* level. This method can lose information leading to performance lower than feature level combination. The potential loss is even greater if the dependence of classifier performance on image quality and pattern are not considered. On the other hand, *decision* level combination promotes information filtering and makes the combination process manageable. In this dissertation we will focus on the combination problem at the decision level rather than at the feature level.

The problem of combining the decisions of multiple classification algorithms has been tackled by many researchers [38], [43], [65], [83]. In general, based on the output provided, there are three types of classifiers:

- (i) Type I, outputs a unique class label indicating the most probable class to which the input pattern belongs
- (ii) Type II, outputs a ranked list of part or all of class labels while a higher a class label is in the list, the more probable it is that the input pattern belongs to the corresponding class
- (iii) Type III, assigns to each class label a measurement value which indicates the degree to which the corresponding class pertains to the input pattern.

Classifier types play an important role in determining the combination method used.

1.1 Objective

Our objective is to develop a general theory of combination that considers the classifier at an abstract level. Information about feature types, training methodology and classifier characteristics can contribute to the design of a good, efficient combinator. On the other hand, not using classifier specific information allows a combination theory to be general.

We study classifier combination under complete abstraction of classifier details. In order to completely separate the design of classifiers and the design of combinators, we will treat each classifier as a black box whose sole interface to the combinator is a set of transformed decision vectors. Since the decision vector transformation is classifier dependent, it has to be provided. Also provided is a set of labeled data whose distribution reflects the environment in which the combined system will operate. All the necessary knowledge about classifiers can be derived through observations of their responses to the observation data. The classes to be recognized in this observation data set are also predetermined. This black box treatment of classifiers can be correlated, trained on different data sets, generating mixed types of outputs, and behaving in various ways. Under these presumptions, we wish to develop a combination theory.

Several issues need to be considered.

1. Combining classifiers is particularly useful in cases when the classifiers are very different in terms of their use of different features sets, training samples, cluster algorithms and measures. The measurement values output by of such classifiers can be interpreted in many different ways: distance, probability, fuzzy measures and so forth. Most known methods of combination simply assume that the scores are probabilities. A solution to this problem calls for a methodology that

will transform the output scores of all kind into probabilities.

2. Some combination methods require classifiers to be *independent*. Since these techniques are simplified based on the independence assumption, they are prone to performance degradation when classifiers are correlated. In reality, classifier independence is difficult to achieve and very hard to detect. In the event that correlation is detected, it is still advantageous to combine them. This calls for the ability to handle correlated classifiers.
3. Practical systems often glue together many modules, where each module has its own output and thresholds that control the accuracy of the module. While each individual module can be optimized by tuning the relevant parameters, it is a non-trivial task to optimize the entire system. Clearly, the optimal setting of parameters for the best performance of the entire system is not simply a choosing of the parameters that optimize the individual modules- for the modules could be competing.

When the number of modules and the parameters are few, manually combining the modules by choosing of all parameters is possible by a "trial and error" mechanism. However, when the modules are many, other methods have to be adopted. Further, most current systems take an approach of "if then else" rules where the parameters are linear and are compared with threshold values. We address these issues in a methodical way.

4. Information about classifier dependencies on image characteristics manifested in classifier output patterns in relation to input images has been largely ignored by most combination algorithms. The fact that classifiers tend to make mistakes on different types of images suggests that combination methods should have the ability to recognize that dependency and dynamically adjust its strategy. The responses of classifiers may be related to each other in an intricate way making

that it is impossible to accurately capture their behavior with a single model. An accurate approximation can be obtained if a separate model is applied to each region of the input space based on some kind of partitioning. Further generalization of this concept would allow a dynamic combination based on a *fuzzy* partitioning of the input space. This generalization is beneficial when definition of mutually exclusive partitions is inappropriate.

1.2 Literature Review

In recent OCR research, combination of several different algorithms has been accepted as an effective method to improve the recognition performance [28], [31], [43], [44], [38], [67], [85]. Many examples of such multiple classifier systems can be found in the literature on handwritten digit recognition [32], [73], [44], [55]. Similar approaches have been used in machine print word recognition [28], cursive word recognition [61] and other pattern recognition problems [21]. Some of the recognition systems use specifically designed combination methods [32], [75], [77], [67], while others use more general combination models [83], [5], [50], [30].

Methods developed for parallel classifier combination include voting principle [83], Bayesian formalism [44], Dempster-Shafer theory [4], neural networks [44], statistical regressions [27] and fuzzy integrals [76]. Each method has its own merits, as well as limitations. We will only briefly explain the fundamental concepts and point out related works.

The voting principle is the simplest method of combination where the top candidate from each classifier constitutes a single vote. The final decisions can be made by the majority rule (over half the votes) [73], [74], plurality (maximum number of

votes)[44], weighted sum or other variants [83]. An interesting theoretical analysis based on the majority rule is provided in [42]. The method of Borda count [5], which sums up the rank values of classifiers, can be considered as a generalization of the voting principle to Type II. Ho investigated its use for word recognizer combination in her work [28]. Friedman's procedure used by Tubbs and Alltop [79] is equivalent to the Borda count. The focus is on calculating the confidence of the combined decision under the normality assumption.

The Bayesian approach seeks to approximate the posterior probability by directly estimating various probability densities [13]. The problem of classifier combination can be viewed as estimating the class posterior probabilities conditioned on classifier decisions. Based on different assumptions and different techniques for probability estimation, variation of the Bayes rules are derived. Simplifications to independent Type I classifiers are most commonly used [83] [44]. Although the decision theory itself does not assume classifiers are independent, this assumption is almost always adopted in practical implementation to reduce the exponential complexity of probability estimation.

A well known theory in evidence combination is the Dempster-Shafer formalism [85] [4] [68]. Its rule of combining individual evidence to form the belief for a single hypothesis reduces to Bayesian inference when all knowledge sources support only simple hypothesis and are completely accurate (no conflicts of beliefs and ignorance) [4]. The difference between the probabilistic approach of Bayesian inference and Dempster-Shafer theory can be found in Barnett [4] and Pearl [60]. Therefore, Dempster-Shafer theory is a more flexible model of reasoning. It takes uncertainty and ignorance into account. Several extensions to the theory have been made. Han and Kashya proposed a modification to improve its robustness when faced with highly conflicting

evidence and a uniform framework to integrate belief combination and belief propagation [24]. Garvey, Lowrance and Fichler proposed a formalism based on *evidential propositional calculus* which represents the belief in a proposition by a subinterval, establishing the *plausibility* and credibility of that proposition [23]. Erkmen and Stephanou incorporated fractal theory into Dempster's theory of reasoning to represent the belief functions, achieving multiple levels of granularity of the frames of discernment [15]. The fundamental property of Dempster-Shafer theory to aggregate beliefs based on uncertainty reasoning lends itself to classifier combination tasks. Mandler and Schurmann [50], Xu, Krzyzak and Suen [83] Lu and Yamaoka [46] applied Dempster-Shafer theory, with simplifications, on handwritten digit and character recognition tasks. Nevertheless, the usefulness of this approach remains limited because Dempster's rule of combination is conditioned on the fact that all bodies of evidences are independent and the complexity of calculating the orthogonal sum is exponential. In addition, the performance of this theory depends on the assignment of basic probability mass, which lacks rigid mathematical formulation. Ling and Rudd generalized the theory to incorporate pairwise dependencies between knowledge sources [47]. However, the degree of dependency needs to be preassigned. Almost all practical implementations of the theory reduce the problem complexity to $O(n)$ by making additional assumptions [4] [83].

Artificial neural networks provide an alternative in modeling the decision process. Although optimality is not guaranteed, good approximations are usually obtained without elaborate hand crafting. Network learning is, for the most part, automatic, and the decision process can be parallelized. It has been shown that a multi-layer perceptron trained by back propagation network [51] finds a least mean square approximation to the Bayesian decisions [45], [64], [80], [33]. The direct estimation of posterior probability endows neural networks with the theoretical support necessary

for their practicality. Their application in classifier combination can be found in a system for postal address interpretation [69]. A comparison of neural networks with other combination methods for handwritten digit recognition is given in [44]. Hunag and Suen describe a neural network combinator that compares favorably to other combination methods [29]. The work of Noumi is another example of a multi-layer perceptron application to classifier combination [56]. Both the neural network and regression approaches, being a direct estimation of the posterior probability, are not dependent on the classifier independence assumption. This is a major advantage over the other methods.

In contrast to the multi-layer perception approach, which is regarded as a form of non-parametric regression, the parametric regression approach estimates the parameters of a predetermined regression model. A popular model used for binary response variables is the logistic function [35]. The efficiency of this method compared to discriminant analysis has been established in [14], [59], [62]. Its first use in word classifier combination was introduced by Ho [26], [27], [28]. A simple regression model based on rank values was shown to be more effective than the highest rank or Borda count method [28]. In addition, the robustness of this approach against correlated classifiers was noted. Other regression models such as the polynomial have also been used in multiple classifier systems [19], [21].

Numerous other approaches have been used to perform combination in particular applications. For example, Takahashi and Griffin use pairwise discrimination to reorder the top decision of a base character recognizer [77]. Sabourin and Nagy devised a K-nearest-neighbor combination for two nearest neighbor classifiers that compares favorably over several other combination methods [67]. A measure based on *relative density* is used by Ho [25] to combine independently trained LVQ classifiers [41]. The combination of pairwise discriminator and multistage design in [18], [20] can also be

considered as examples of classifier combination. Fuzzy integration of classifier reliability and objective evidence has also been successfully used for classifier combination [76], [84].

The majority of work in this field has been on Type I classifiers because the number of possible combinations are discrete and finite. Xu Krzyzak and Suen considered Type I classifier combination by voting principle and Dempster-Shafer theory [83]. Huang and Suen take the Bayesian approach to build a theoretically optimal framework for combining type I classifiers [30].

In the combination of Type II classifiers, although the number possible combinations are still discrete and finite, they suffer from combinatorial explosion. Ho [28], [25] used logistic regression under the assumption that each classifier recognized all classes equally well. For large number of classes, the method of union and intersection can be effective in reducing the number of candidate choices [26]. Tubbs and Alltop used a non-parametric procedure to combine ranked lists [79]. Under the uniform distribution assumption, their procedure reduces to summation of ranking values. The reason for making the uniform distribution assumption is not to simplify the combination rule as much as to force the summed value to an asymptotic chi-square distribution so that a probabilistic measure of confidence can be calculated.

Combining Type III classifiers is the most difficult because of the potentially infinitely many possible combinations of confidence values and the complicated relations between class confidences to classifier performances. The problem is further compounded by the fact that not all the classifier scores are on a uniform scale. Several researchers have tackled this problem. Mandler and Schurman described a method for transforming distance measures to confidence values suitable for basic probability assignment in Dempster-Shafer theory based on inter-class and intra-class

distributions [50]. In the work of Xu, Krzyzak and Suen, simple averaging is used for combination [83]. The averaging method clearly did not take difference in classifier performance, score scales and classifier dependence into consideration. The technique of fuzzy integral described in [76] combines classifiers taking both classifier reliability and decision confidence into account. This method attempts to resolve the difference between classifier performance and observed evidence for a given sample by applying a min-max rule. However, the scores are again put on the same scale as probabilities and compared.

Combination of mixed types of classifiers has not been explicitly addressed in the past. Obviously, algorithms which combine only Type I classifiers don't have such concerns. Most of the methods can be modified to work with Type I and Type II mixtures. However, generalization to include continuous data of Type III and finite valued Type I and Type II data requires special treatment. Therefore, type transition becomes non-transparent to the combinator. In other words, a different combinator is needed for mixed types of classifiers. An easy solution to achieve type homogeneity is by reducing information in all classifier outputs to the lowest level. However, in this particular case the loss of information affects combination performance. Although many researchers acknowledge the value of additional information contained in higher types [28], [30], there is no study comparing the contributions of each level of information for a single classifier.

The effect of correlated classifiers on combination performances has not gone unnoticed. The assumption of *classifier independence* is often alluded to in the literature [83], [44], [50]. In some situations, the assumption is forced by the combinatoric nature of density estimation in high dimensionality, as is the case in the Bayesian approach proposed by Lee [44]. In others, the independence assumption is inherent to the theory itself, as in the case of Dempster-Shafer theory and voting principle[83],

[50]. Huang and Suen address this problem by deriving a Bayes optimal combinator for Type I classifiers and suggest a practical solution to overcome combinatoric explosion [30]. However, the method is not generalizable to Type II and Type III outputs and the problem complexity is still exponential. Ling and Rudd proposed an extension to the Dempster-Shafer theory to account for pair-wise classifier dependencies [47]. Nonetheless, the degree of dependency needs to be preassigned, a difficult task in practice. The neural network and regression approach use direct function approximation to the posterior probabilities and, therefore, are immune to correlated classifiers by nature of their solution finding strategy. Consequently, all combination methods utilizing these approaches are not be affected by correlated classifiers. Lee [44], Huang [29] and Noumi [56] all used a single hidden layer multi-layer perceptron for combination. Examples for the regression approach include work of Ho [28] and Franke [19], [21].

The idea of redundant classifier elimination is similar to variable [36], [78]. Research in this area in the context of classifier combination is limited. Ho suggested removal of classifiers whose decisions are always reinforcing that of others for the highest rank and Borda count method. In the case of the logistic model which has one weight per classifier, the decision for elimination can be based on the weight magnitude [28]. The major reason for the anonymity of this topic is the small number of classifiers involved and that the best combination can be obtained by hand crafting (Chou, etc. [11]).

Estimating classifier performance dependencies on images is closely related to one type of dynamic classifier selection. The difference is in the selection criteria. Several researchers have successfully explored this approach [28], [67]. One method suggested is to measure classifier performance based on mutually exclusive partitioning of the data. It is shown that selection based on whether top choices of classifiers

agree improves the performance of the combinator [28]. Sabourin, etc. base their dynamic selection on parameters which have “the highest mutual information with correctness” [67]. The application of information theory in this task is very interesting. However, their selection criteria are based on *post* classifier decision factors and can not capture classifier dependency on image characteristics. Preconditioned dynamic feature/classifier selection is suggested by Ho [25]. However, no distinction between the two types is made and no concrete algorithm is given. Nowlan’s competing network system addresses such dependencies by training a gain control network to arbitrate subtask learning in a multiple expert environment [57], [58]. Through *soft* competitive learning, experts are trained to specialize in subtasks so that their joint effect, controlled by a gating network, achieve overall optimization. The *preconditioned* gating network is trained on input or *derived input* in conjunction with the learning experts. Another example of preconditioned selection appears in the work of Xu and Jordan [82]. In their posterior probability estimation using a mixture of Gaussian model, the parameters for the mixtures are explicitly expressed as functions of the input image, suggesting they can be dynamically modified. The parameters are calculated using EM algorithm [52].

1.3 Outline of the Dissertation

Following is the outline of this dissertation. Our focus will be on two broad areas.

- Measures to normalize scores output by classifiers and
- Applications

1.3.1 Measures

1. Mapping Scores to Probabilities:

We will develop the DPPS (Deriving Pseudo-Probabilities given Scores) model for expert systems. In order to completely separate the design of the expert system from the methodology, we will treat the expert system as a "black box", that transforms the input into a decision vector. We will also get an uniform measure as a powerful tool for evaluation of performance of different expert systems, as well as the interpretation of scores.

2. Fuzzy Measure:

We will use a fuzzy measure based on the topology of the "blind" set of training samples and pre-determined neighborhoods to construct the fuzzy membership function. This will allow a non-linear partitioning of the space, formed by scores that the expert returns as outputs. As a result of this new nonlinear partition of the score space we will improve the performance of a given expert and get a new uniform measure that can be used in combination of different types of experts.

3. Accuracy Maps:

Recognizers usually differ in their confidence measures. A judicious combination of these measures can be made only when they are on the same scale. We will develop a way to transform confidence measures from different recognizers to a "special" confidence, we call accuracy. For a given set of "large enough" as well as "properly representative" training samples, the accuracy measure can be obtained by transforming the confidence using sufficient training samples. Usually, the accuracy is an increasing function of confidence.

Let us assume that the training set is adequately representative. We will create the accuracy approximation by utilizing the training samples. We assume that

the confidence measures of all classifiers are continuous, that is, the confidence values can be any point in $[0, 1]$. We will introduce a simple and efficient method to implement such a transform.

1.3.2 Applications

1. Confidence Combination in Multi-expert Systems:

Previous studies have developed many different approaches for expert combination. For experts of Type I for which only labels are available, voting algorithms are used. Label rankings are used with experts of Type II. In case of Type III experts with measurement values interpreted as posteriori probabilities, a Bayesian technique is often applied for experts assemblies. If the expert output is interpreted as fuzzy membership or evidence values, fuzzy rules and Dempster-Shafer approach are used. Also there are cases of expert combination, where the output of the expert is used as a new feature and a new expert (neural network) is built to perform the combination.

We will focus on a simplified version of Type III multi-expert combination. Each expert uses its own representation, i.e. measurements extracted from the input pattern are unique to each expert. We will create a new accuracy measurement scale, uniform for all experts and use it as a basis for the expert combination. We will develop an optimal combination scheme which requires extremely large amount of training data as well as memory space. So we introduce an empirical scheme to approximate the optimal scheme so that the requirements on the amount of training data and memory is practical. We first develop the accuracy measure which is common for all experts. Next, we characterize each expert with a family of accuracy maps. Next, we build accuracy combination maps with a special synthetic function. Finally, we construct synthetic accuracy maps

for the combined confidences. We also propose a simplified combination scheme which requires less training data while sacrificing some accuracy. We finally discuss the optimal rejection threshold for the final recognition decision.

2. Handwritten Address Interpretation Systems:

We propose a new Fuzzy set technology in the application of a Handwritten Address Interpretation System. The modules that constitute the system are: (i) Address Block Location, (ii) Line separation module, (iii) Word Separation Module, (iv) ZIP Code recognition module, (v) Street number recognition module, and (vi) Street name recognition module. Each of these modules have submodules and face the same problems at a micro-level. The ZIP Code recognition and Street number modules invoke a numeral recognition module and a Street number recognition module invokes a character recognition module. Each of the modules and sub-modules are controlled by many parameters.

In the final analysis, a particular recognition of the address block is deemed either worthy of being accepted as correct or rejected as unreliable. The idea is to accept as high a percentage of addresses as possible while making a minimum of errors. A conservative setting of the thresholds at the final stage would cause many correct results to be rejected in order to keep the error rate low.

3. Architecture for Serial Classifier Combination Using Entropy Measures:

Classifiers can combine in either series or in parallel. We propose to develop a general theory for serial combination of classifiers. We will discuss the theoretical and complexity issues pertaining to lexicon reduction and the serial combination of classifiers.

The method of parallel combination of classifiers would submit the same lexicon to both classifiers and combine the results using a variety of methods, such as logistic regression and Borda Count. On the other hand, in serial combination,

classifiers that operate later in the engine deal with smaller lexicons. In fact, lexicon reduction is central to the serial combination methods.

The motivation stems from the difficulty confronted by a designer in choosing the correct architecture. Empirical methods are typically used by researchers. However, as the number of classifiers increases and the size of the lexicon increases, the situation quickly gets out of hand. Our objective is to develop the theory that will provide the guidelines for choosing the best method of serial combination.

Chapter 2

Measures

Pattern recognition has been an active area of research for many years and fostered the development of many recognition techniques. Combining the results of these methods can improve the overall performance. Combining classifiers is particularly useful in the cases when the classifiers are very different (i.e use different features sets, training samples, cluster algorithms, measures, etc.). The outputs of such classifiers can be interpreted in many different ways: distance, probability, fuzzy measures and so forth. In addition, most of the known methods of combination assume that the scores returned by classifier are probabilities. A solution to this problems is to create a methodology that will transform the output scores of classifiers into probabilities. This chapter is an attempt to address these needs.

2.1 Deriving Pseudo-Probabilities of Correctness Given Scores (DPPS)

Our main goal is to describe how to apply the DPPS (Deriving Pseudo-Probabilities of Correctness Given Scores) model to recognition systems. Although a thorough theoretical investigation is beyond the scope of this section, we will demonstrate the effectiveness of the method by experimental results. The benefits of the application of DPPS is two-fold:

- (i) improvement of the pattern recognition and
- (ii) further applications of the derived pseudo-probabilities.

Experiments show an improvement in the range 0.56-1.5% in the decision making of a classifier when using the new pseudo-probabilities. We also get a uniform measure as a powerful tool for evaluating the performance of different recognition systems and interpreting their scores.

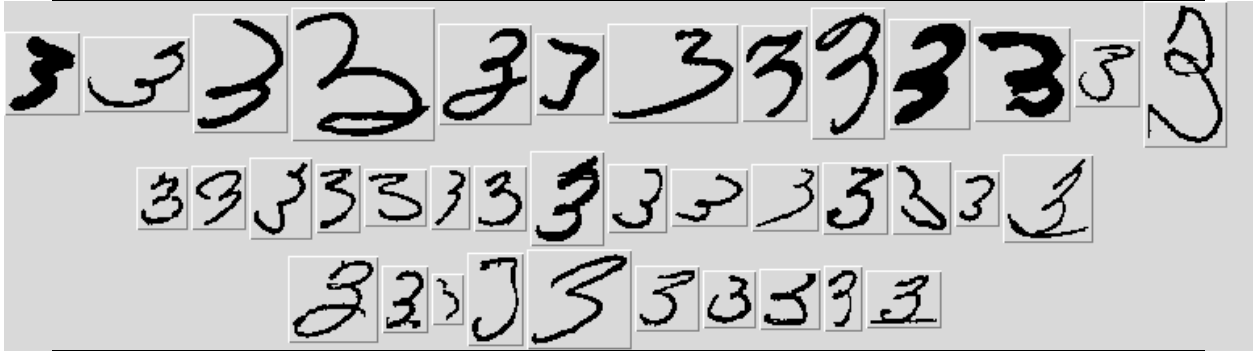


Figure 2.1.1: The Partition Set V_3^3 from GSC

A classifier of type III takes an unknown pattern for an input and outputs a decision vector of numerical scores, that classifier assigns to each of the classes. We will describe how DPPS transforms the output of the classifier into pseudo-probabilities of correctness, associated with each class. In order to completely separate the design

of the classifier from the DPPS methodology, we will treat each classifier as a "black box", that transforms the input into a decision vector.

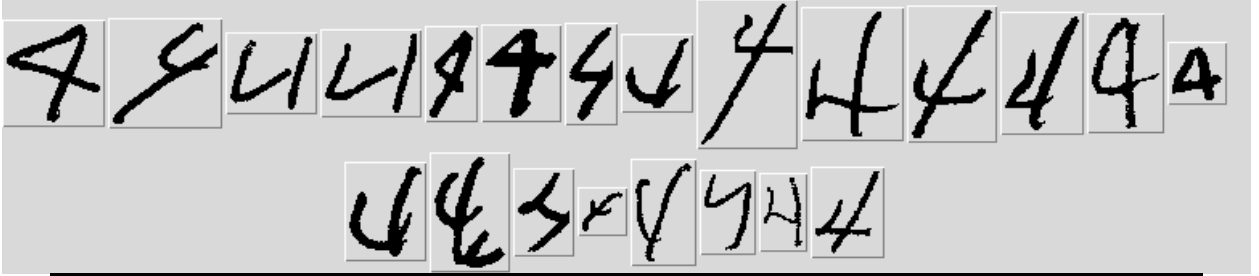


Figure 2.1.2: The Partition Set V_4^1 CP classifier

We will support the DPPS concepts by experiments using the following classifiers: **Polynomial** classifier ([71]), **Gradient** classifier ([70]), **GSC** classifier using k-nearest neighbor ([17]) and **CP** (a logistic regression combination of gradient and polynomial classifiers), developed at CEDAR. The '**Polynomial**' classifier is based on a polynomial discriminant algorithm trained to extract a relative weighting for each feature in each class. The **Gradient** classifier encodes local contour variations of a character image into a binary feature vector. **GSC** classifier uses 512 symbolic multi-resolutional features, which measure the image characteristics at local, intermediate, and global scales. The 192 **Gradient** features measure edge curvature in a neighborhood of a pixel and provide information about stroke shape on a small scale. The 192 intermediate **Structural** features measure short stroke types which span several pixels and give useful information about stroke trajectories. The 128 **Concavity** features are used to detect stroke relationships which can span across the image. Features at all the 3 levels, **G**, **S**, and **C** are combined in a 512 binary feature vector and a weighted k-nearest neighbor algorithm is used for classification.

2.1.1 An Intuitive Problem Description

Following is the analogy that should clarify how the DPPS algorithm works. Teacher A and Teacher B want to evaluate the proficiency of students in Math. Both teachers give a question to n students (S_1, \dots, S_n) and grade the students' responses. Teacher A examines the student's response and gives a score of 80% to the student (S_j) with the best answers (other students get lower scores). Teacher B gives a score of 90% to the same student (S_j) which is also the highest score in his grading. Now observing the above, the following questions can be asked:

1. Is student S_j the most proficient in Math among the students examined?
2. Is the opinion of Teacher B about S_j 's proficiency stronger than that of Teacher A, given that Teacher B gave the student a higher score?
3. When Teacher A (B) give the best student a score of 80% (90%), are they correct in their selection of the best student?

An intuitive answer to the first question would be that indeed student S_j is the most proficient among the students $S_1 \dots S_n$, if he scores the highest consistently over many exams. This eliminates the possibility of chance occurrence. However, it must be ensured that all the exams administered are of equal difficulty.

An intuitive answer to the second question is that the information given is insufficient to make the determination. The grading policies of both teachers have to be studied over a large number of tests to qualify their grading behavior. In one teacher's mind 80% can mean more than what 90% means to the other teacher. The reason that comparison across teachers is difficult is because each teacher has a different notion of what a particular score means, i.e score are on different scales.

It is the third question we address here – in this case the set of all classes are the students, the classifier plays the teacher role and the input pattern is in the role of an exam. Our objective is to assess to what degree the classifier is correct in its ability to label unknown patterns and in particular, its ability to choose the best class. The degree of correctness over a large number of trials should provide the probability of correctness.

It is reasonable to assume that any class ω has the same chance of being the truth, i.e all classes are distributed uniformly. We circumvent the difficulty of computing $P(\omega|x)$ by falling back to the classic definition of probability. According to the definition probability is the frequency of outputs of an event of interest over many trials. However, given a pattern, a recognizer provides the same score for each class no matter how many times the recognition is performed. In other words, there is no randomness inherent in this process. So, one way to introduce randomness in this case is follows: instead of recognizing over and over the same pattern with the same result we generate many recognition trials with patterns “similar” to the input pattern. A set of “similar” patterns is generated in the re-training process by collecting the input pattern and its neighbors and then it is used to simulated the repeated recognition trials. In our analogy, this is the process of creating more exams for the students, where all the exams are of the same difficult level as the first one. Let us note that the “retraining” phase brings more knowledge and potential improved recognition rate.

Given an unknown input pattern X , we are interested in finding what is the chance (probabilistically) that a given class will be among the top choices (confusion set).

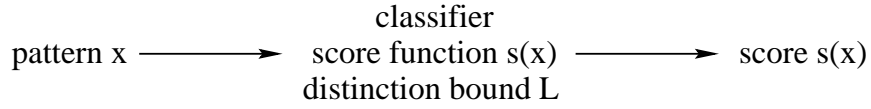


Figure 2.1.3: A model of an one dimensional classifier: if $s(x) < L$, then x belongs to class $\omega_1 \in N(a, \sigma^2)$, if $s(x) > L$, then x belongs to class $\omega_2 \in N(b, \sigma^2)$

2.1.2 Theoretical Foundation of DPPS Methodology

Let us consider the following hypothetical example [13]. Suppose that a fish packing plant wants to automate the process of sorting incoming fish on a conveyor belt according to the species. As a pilot project it is decided to try to separate sea bass from salmon using optical sensing. A camera is set and it takes some sample images. The camera's signals are processed and a measurement value is returned, based on the classification made (i.e the returned score is compared to a threshold L).

This example describe a one dimensional classifier (see Figure 2.1.3). Let us assume that class ω_1 of salmon is normally distributed with mean a and deviation σ with respect to the scoring function. Let us also assume that class ω_2 of sea bass is normally distributed with mean b and deviation σ with respect to the scoring function. Also, without loss of generality, we can assume that $a < L < b$ (see Figure 2.1.4).

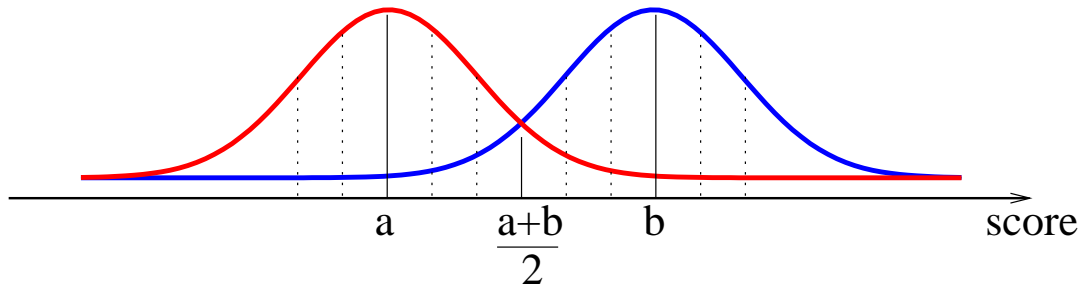


Figure 2.1.4: Distribution of salmon and sea bass on the scale of scoring function

Theory: The classifier is theoretically optimal when $L = \frac{a+b}{2}$

Proof: see [13].

Suppose that our classifier is not optimal, i.e. $L \neq \frac{a+b}{2}$. Our goal is to show that applying our methodology of computing pseudo-probabilities of correctness, we modify the classifier into the optimal one, i.e given unknown pattern x , if $s(x) < \frac{a+b}{2}$, then $x \in \omega_1$ and if $s(x) > \frac{a+b}{2}$, then $x \in \omega_2$.

Let V be the validation set with $2n$ samples, n samples from the first class ω_1 and other n samples from the second class ω_2 . We will assume that n is large enough, so the set V represents both classes well.

Let

$$f(z, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(z-\mu)^2}{2\sigma^2}}$$

be the density function of a normal distribution with mean μ and deviation σ .

Let us consider the partition of the real line (\mathbb{R}) using the following points

$$\begin{aligned} \{\dots, y_{-3} = 2a - b, y_{-2} = \frac{3a - b}{2}, y_{-1} = a, y_0 = \frac{a + b}{2}, \\ y_1 = b, y_2 = \frac{3b - a}{2}, y_3 = 2b - a, \dots\} \end{aligned}$$

Let $X_{s(x)}$ denote the confusion set of an unknown pattern x , i.e the set of all samples from the validation set V that are similar to pattern x in view of the scoring function. More precisely, if $y_k < s(x) < y_{k+1}$ for the points y_k and y_{k+1} from the partition, then

$$X_{s(x)} = \{v \in V : y_k < s(v) < y_{k+1}\}.$$

Let N_i is the number of the images in $X_{s(x)}$, which belong to the class ω_i , $i = 1, 2$. Then,

$$N_1 \approx n \int_{y_k}^{y_{k+1}} f(z, a, \sigma) dz$$

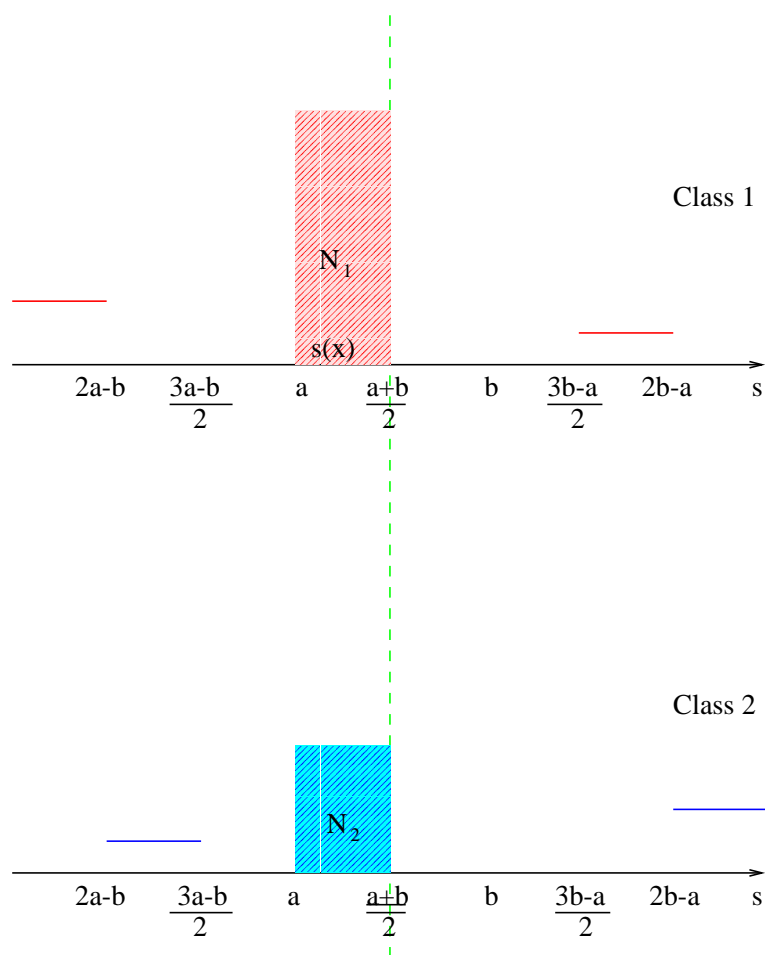


Figure 2.1.5:

$$N_2 \approx n \int_{y_k}^{y_{k+1}} f(z, b, \sigma) dz.$$

Lets us make the following observation. Since $f(z, a, \sigma) > f(z, b, \sigma)$ for $z < \frac{a+b}{2}$ and $f(z, a, \sigma) < f(z, b, \sigma)$ for $z > \frac{a+b}{2}$, if $s(x) < \frac{a+b}{2}$, then it is obvious that the number N_1 of samples of class ω_1 in the confusion set $X_{s(x)}$ is greater than the number N_2 of samples of class ω_2 in the set $X_{s(x)}$ (see Figure 2.1.5). Then, the pseudo-probabilities of correctness are computed as follows:

$$\xi_i(x) = \frac{N_i}{|X_{s(x)}|} \quad i = 1, 2$$

implies that $\xi_1(x) > \xi_2(x)$ if $s(x) < \frac{a+b}{2}$ and $\xi_1(x) < \xi_2(x)$ if $s(x) > \frac{a+b}{2}$.

The same idea can be applied to the following model of a classifier (Figure 2.1.6)

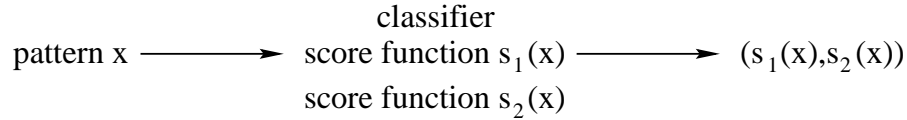


Figure 2.1.6: A model of a two dimensional classifier: if $s_2(x) < s_1(x)$, then x belongs to class ω_1 , otherwise x belongs to class ω_2

in order to improve the classifier performance.

2.1.3 DPPS Algorithm

For simplicity, we assume that our classifier is the **GSC** digit recognizer [17]. It takes a digit image as an input and return a decision vector of the following format

$$(c_0 \ s_0 \ c_1 \ s_1 \ c_2 \ s_2 \ c_3 \ s_3 \ c_4 \ s_4),$$

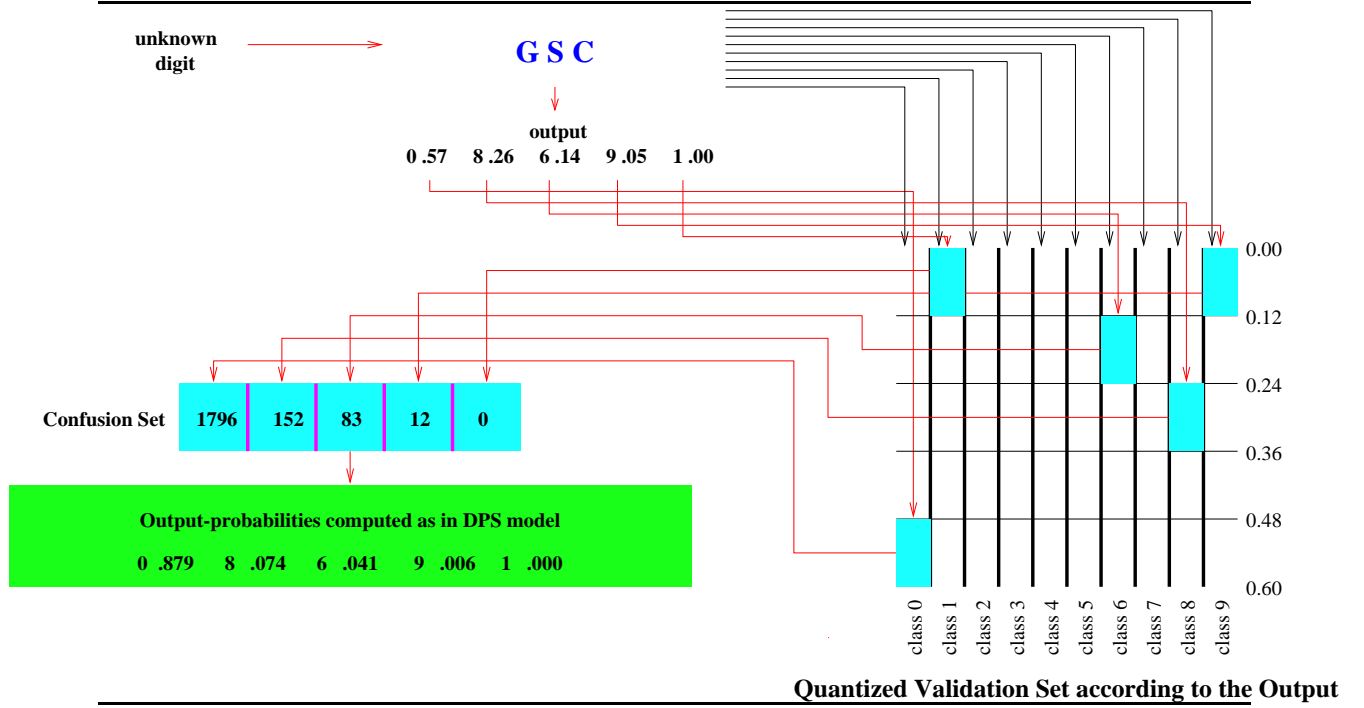


Figure 2.1.7: The DPPS Algorithm in this application

where c_0 is the first choice class with the highest score s_0 , c_1 is the first choice class with highest score s_1 and so on. All scores are real numbers between 0 and 0.60.

We assume that all classes have the same chance of being the correct one. We use the set V of 20000 digit images as a validation set, that we use for re-training of the classifier. Each class is represent in V by 2000 digits. We propose four different techniques to quantize the validation set. They are described as follows:

1. Noisy class histogram (input partition)

In this case we quantize the set V with respect to the input. Then $V = \bigcup_{i=0}^{i=9} V_i$, where V_i is the subset of V containing all digits in V that belong to class ω_i . We divide each of the sets V_0, V_1, \dots, V_9 into q subsets $V_i = \bigcup_{j=0}^{j=q-1} V_i^j$ according to the

score assigned by the recognizer to the top choice class. The set V_i^j contains all digits of class ω_i whose first choice score is in the interval $[0.60j/q, 0.60(j+1)/q)$.

2. Pure class histogram (input-output partition)

In this case we quantize the validation set with respect to the input and output. We proceed in the same way as in the case of noise class histogram and then eliminate from the partition V_i^j all digits which are not recognized correctly, i.e. the first choice class is different from the class the digit belongs to.

3. Mixed class histogram (output partition)

In this case we quantize the validation set with respect to the output. Then $V = \bigcup_{i=0}^{i=9} V_i$, where V_i is the subset of V containing all digits in V that have class i as a first choice class in the decision vector, returned by the recognizer. We proceed in the same way as in the case of noise class histogram.

4. d-neighborhood

For a given input digit image we use the entire validation set to form an appropriate d-neighborhood that contain all digits of the validation set that are "similar" to the input in distance topology, instead of using the one generated by the fixed quantization of the validation set. In other words, for given d , the d-neighborhood contains all samples from the validation set V , for which the following statement is true: The first choice class is among the classes in the decision vector and the difference of the corresponding scores is less than d . We will call that d-neighborhood a confusion set for the input. We also consider two possibilities-pure or mixed d-neighborhoods, which are described in the same way as in the case of pure and mixed histogram.

Let us note that all digits in each of the partition set V_i^j , $i, j = 0, 1 \dots, 9$ are "similar"

method	test set 3-30398 digits	
	correctly recognized	% of correctly recognized
GSC	29532	97.15%
Noise	29548	97.20%
Pure	29583	97.32%
Mixed	29632	97.49%
Knn	29701	97.71%

Table 2.1.1: Comparison of GSC and DPPS using different techniques

from the view point of the recognizer (Figure 2.1.1 and Figure 2.1.2).

Once we quantize the validation set, our next goal is to form the confusion set X_a for a given pattern a , that contains all "similar" samples from the viewpoint of the recognizer. In this particular application we have chosen $q = 5$, i.e. we have divided each of the sets V_0, V_1, \dots, V_9 into 5 parts. In case of noisy class histogram, V_i^j contains all digits in the validation set V of class ω_i , such that their top score in the decision vector, produced by the recognizer, is in the interval $[0.12j, 0.12(j+1))$. In case of pure class histogram, V_i^j contains all digits in the validation set V of class i , such that their first choice class is ω_i and the score in the decision vector, produced by the recognizer, associated with that class is in the interval $[0.12j, 0.12(j+1))$. In case of mixed class histogram, V_i^j contains all digits in the validation set V , such that their first choice class is ω_i and the score in the decision vector, produced by the recognizer, associated with that class is in the interval $[0.12j, 0.12(j+1))$.

In the first three cases (noise, pure and mixed class histograms) we proceed as follows. Let

$$(c_0 \ s_0 \ c_1 \ s_1 \ c_2 \ s_2 \ c_3 \ s_3 \ c_4 \ s_4),$$

be decision vector associated with the digit a , with $0.12k_j \leq s_j < 0.12(k_j + 1)$ for

class	size of the test set	GSC		DPPS	
		correctly recognized	% of correctly recognized	correctly recognized	% of correctly recognized
0	6986	6833	97.81	6852	98.08
1	6554	6496	99.12	6513	99.37
2	3954	3759	95.07	3764	95.19
3	2238	2170	96.96	2177	97.27
4	3191	3070	96.21	3075	96.36
5	1556	1497	96.21	1502	96.53
6	1474	1442	97.83	1446	98.10
7	1445	1391	96.26	1405	97.23
8	2049	1945	94.93	1958	95.58
9	951	929	97.69	943	99.15
total	30398	29532	97.15	29635	97.49

Table 2.1.2: Performance of GSC and mixed class histogram DPPS on different classes

class	size of the test set	GSC		DPPS	
		correctly recognized	% of correctly recognized	correctly recognized	% of correctly recognized
0	6986	6833	97.81	6867	98.30
1	6554	6496	99.12	6518	99.45
2	3954	3759	95.07	3770	95.35
3	2238	2170	96.96	2183	97.54
4	3191	3070	96.21	3082	96.58
5	1556	1497	96.21	1509	96.98
6	1474	1442	97.83	1453	98.58
7	1445	1391	96.26	1411	97.65
8	2049	1945	94.93	1963	95.80
9	951	929	97.69	945	99.36
total	30398	29532	97.15	29701	97.71

Table 2.1.3: Performance of GSC and d-neighborhood DPPS on different classes

class	size of the test set	Polynomial		DPPS	
		correctly recognized	% of correctly recognized	correctly recognized	% of correctly recognized
0	6986	6586	94.27	6613	94.66
1	6554	5841	89.12	5858	89.38
2	3954	3484	88.11	3538	89.50
3	2238	1858	83.02	1889	84.41
4	3191	3011	94.36	3034	95.08
5	1556	1428	91.77	1536	98.71
6	1474	1352	91.72	1382	93.76
7	1445	1344	93.01	1379	93.55
8	2049	1650	80.53	1768	86.29
9	951	807	84.86	829	87.17
total	30398	27361	90.00	27826	91.53

Table 2.1.4: Performance of Polynomial classifier and DPPS on different classes

$j = 0, 1, 2, 3, 4$. Then $X_a = \bigcup_{j=0}^{j=4} V_{c_j}^{k_j}$. For the purpose of better illustrating how we construct the confusion set for the digit a , we will consider the following example. The decision vector produced by the **GSC** recognizer on a unknown digit is

$$(8 \ 0.289 \ 3 \ 0.083 \ 0 \ 0.078 \ 5 \ 0.039 \ 2 \ 0.005).$$

Then the confusion set

$$X = V_8^2 \bigcup V_3^0 \bigcup V_0^0 \bigcup V_5^0 \bigcup V_2^0.$$

In the case of d-neighborhood, we define the confusion set precisely: $X_a = \bigcup_{i=0}^{i=4} W_i^a$, where W_i^a contains all digits for which at least one of the following 5 statements is satisfied:

- the first choice class in the decision vector is c_0 and top choice score is in the interval $[s_0 - 0.06, s_0 + 0.06]$

class	size of the test set	Gradient		DPPS	
		correctly recognized	% of correctly recognized	correctly recognized	% of correctly recognized
0	6986	6597	94.43	6626	94.85
1	6554	6107	93.18	6212	94.78
2	3954	3592	90.84	3613	91.38
3	2238	1989	88.87	2005	89.59
4	3191	3028	94.89	3030	94.95
5	1556	1441	92.61	1445	92.87
6	1474	1374	93.22	1378	93.49
7	1445	1299	89.90	1301	90.03
8	2049	1711	83.50	1824	89.02
9	951	863	90.75	895	94.11
total	30398	28001	92.11	28329	93.19

Table 2.1.5: Performance of Gradient classifier and DPPS on different classes

- the first choice class in the decision vector is c_1 and top choice score is in the interval $[s_1 - 0.06, s_1 + 0.06]$
- the first choice class in the decision vector is c_2 and top choice score is in the interval $[s_2 - 0.06, s_2 + 0.06]$
- the first choice class in the decision vector is c_3 and top choice score is in the interval $[s_3 - 0.06, s_3 + 0.06]$
- the first choice class in the decision vector is c_4 and top choice score is in the interval $[s_4 - 0.06, s_4 + 0.06]$.

Here, $(c_0 s_0 c_1 s_1 c_2 s_2 c_3 s_3 c_4 s_4)$, is the decision vector associated with the digit a .

Once we form the confusion set, then deriving the probability measure is very simple. If $|S|$ denotes the number of the elements (cardinality) in the set S , then the probability the unknown digit a being of class ω_i will be equal to the proportion of the number of elements of the set V_i , presented in the confusion set over total number of

class	size of the test set	CP classifier		DPPS	
		correctly recognized	% of correctly recognized	correctly recognized	% of correctly recognized
0	6986	6612	94.65	6672	95.51
1	6554	6372	97.22	6415	97.88
2	3954	3636	91.70	3661	92.59
3	2238	2094	93.57	2112	94.37
4	3191	3047	95.49	3088	96.77
5	1556	1467	94.28	1483	95.31
6	1474	1398	94.88	1421	96.40
7	1445	1301	90.03	1319	91.28
8	2049	1930	94.19	1947	95.02
9	951	909	95.58	918	96.53
total	30398	28766	94.63	29036	95.52

Table 2.1.6: Performance of CP classifier and DPPS on different classes

elements in the confusion set, i.e.

$$P(a \in \omega_i | V) = \frac{|V_i^{k_i}|}{|X_a|},$$

for $i = 0, 1, 2, 3, 4$

2.1.4 Experimental Results

In this section we give the experimental results from the performance of DPPS methodology in combinations with different classifiers. Table 2.1.1 contains the results from the DPPS method based on the scores given by the **GSC** recognizer. We have conducted experiments, involving all possible techniques described in previous section. Table 2.1.2 and Table 2.1.3 give more detailed comparison of the performances of **GSC** and DPPS model. Table 2.1.4 compares the performance of Polynomial classifier and DPPS model. Table 2.1.5 shows how DPPS compares to the gradient classifier and Table 2.1.6 shows the performance of enhanced CP classifier.

2.1.5 Some Final Remarks and Observations

First, we can see performance of the enhanced by the DPPS methodology classifier. The results also show that the extent of improvement depends on the distribution of the different classes.

Second, if the original recognition system performs more “poorly”, DPPS improves the performance of the system by a greater degree.

Our hypothesis is that the performance of the algorithm will improve significantly provided the training set captures the properties of the recognition system, as well as optimal partition limits, needed for the quantization of the validation set. In some cases, it can be proved theoretically that proper validation set in combination with proper partition can improve the recognition system to the theoretically optimal.

DPPS modified classifier is almost as time consuming as the classifier itself in cases when d-neighborhood approach is not used. In the cases, when we choose to use d-neighborhood DPPS, the better correctness rate is a trade-off for increased classification time.

We have presented an algorithm for deriving probability values given recognition scores assigned to digits by a recognizer of type III. We have argued that these probability values are more informative and useful in a variety of applications including, but not limited to, multiple classifier methodologies.

2.2 Fuzzy Measure and Potential Improvement of Classifier Accuracy

Typical digit recognizers classify an unknown digit pattern by computing its distance from the cluster centers in a feature space. The K-Nearest Neighbor (KNN) Rule assigns the unknown pattern to the class belonging to the majority of its K neighbors. These and other traditional methods adopt a uniform rule irrespective of the “difficulty” of the unknown pattern.

In this section, we propose a methodology which has many salient aspects. First, the classification rule is dependent on the “difficulty” of the unknown sample. Samples “far” from the center which tend to fall on the boundaries of classes are error-prone and hence “difficult”. A “overlapping zone” is defined in the feature space to identify such difficult samples. A table is pre-computed to facilitate an efficient lookup of the class corresponding to all the points in the overlapping zone. The lookup function itself is defined by a modification of the KNN rule.

By reassigning the error-prone samples, we actually redefine the class boundaries. A characteristic function defining the new boundaries is computed using the topology of the set of samples in the overlapping zones.

Our two-pronged approach uses different classification schemes with the “difficult” and “easy” samples. This could enable potential improvement of a recognition system. We have tested this methodology on a large set (30,398) of handwritten digit images. The method described in this section has improved the performance of the GSC digit recognizer described in [17]. The GSC is among the best digit recognizers described in the literature. The method described in this section successfully reduced its error

rate from 2.85% to 1.96%, i.e by 0.89%, which is more than 30% of the initial error. We have tested our method on other available classifiers and have obtained similar results.

Pattern classification consists of assigning a class label to a set of unclassified patterns. Classification typically involves two steps: first, the system is trained on a set of data (already labeled prototype patterns) and then this data set is used to classify a new set of unclassified patterns. When all possible classes are known in advance and the system is trained on a set of already classified patterns, the problem is known as supervised classification. Formally, given the pattern x (taken from the space of all unknown patterns X), a set of c class labels ($\omega_1 \dots \omega_c$) and a set of labeled prototypes, the classifier assigns a label to x such that there is a minimum probability of misclassification. For the purpose of this section we will confine our discussions to the handwritten digit recognition problem where the number of classes is 10.

A typical digit recognizer returns scores corresponding to each class. The scores reflect the degree to which the unknown pattern x represents the prototype(s) of the class. The scores can be represented by the vector $\Phi(x)$, where $\Phi(x) = [\phi_0(x), \phi_1(x), \dots, \phi_9(x)]^T$. Each $\phi_i(x)$ corresponds to class ω_i . $\phi_i : X \rightarrow [0, 1]$ is a function that maps each unknown pattern to a real number between 0 and 1. The performance of the digit recognizer is measured by relative frequency with which it assigns unknown patterns to their correct classes. In other words, the accuracy is the ratio of the number of correctly classified patterns and the number of all classified patterns. Higher the accuracy, more desirable is the recognizer in practical applications.

2.2.1 Issue with traditional techniques

Traditional classification techniques behave as follows. Given an input pattern x , a classifier typically computes the distance of x from the centers of each of the prototype classes in the feature space using some distance metric. The classes are then ranked in increasing order of the distances. In Figure 2.2.1, x is “closest” to the center of the class ω_5 with a distance of s_5 . The 2nd and 3rd classes are ω_6 and ω_0 at distances s_6 and s_0 respectively. It is to be noted that usually $\phi_i(x)$ is inversely proportional to s_i .

It is possible that the test pattern x falls within the boundary of one class but its distance is closest to the center of a neighboring class. Since x , in the example, falls close to the boundary of classes ω_5 - ω_0 and ω_5 - ω_6 there is more than an average chance of misclassification of x .

Figure 2.2.1 illustrates the class boundaries among digit prototypes in a hypothetical feature space. We will assume for the purpose of this discussion that the classes are linearly separable. Samples that are “close” to the centroid of the classes are “stronger” representatives of the class while samples close to the boundaries are poorer representatives and at the same time prone to confusion with the neighboring classes.

So, it is natural for one to consider another level on top of the classifier in order to deal with the patterns that are more prone to misclassification. There are many different techniques that are applicable to this problem.

Using the KNN rule (see [12], [13]) is one way to address the problem stated above. Instead of simply relying on the distance of a sample from cluster centers, one labels x as the majority class in the set of its k nearest neighbors in feature space.

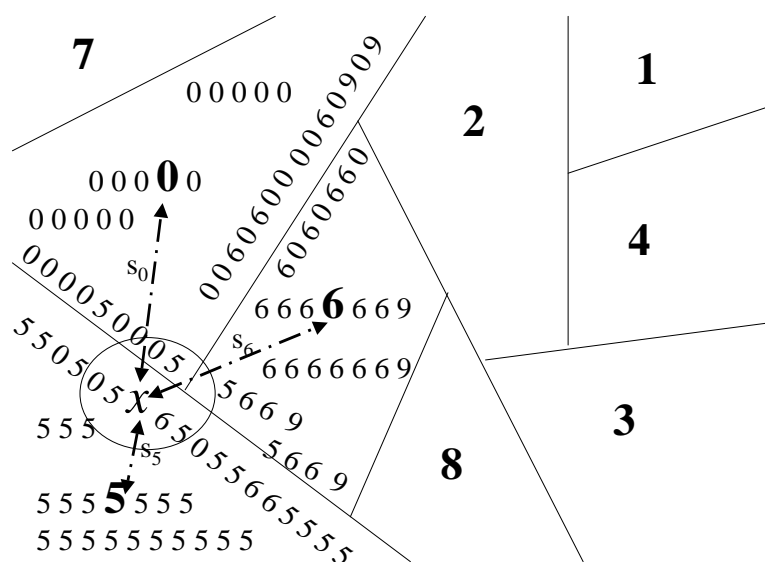


Figure 2.2.1: Feature space of a typical digit recognizer. Samples that are close to the center of a class are strong representatives of the class and are usually correctly classified. Samples that are away from the center (close to the boundaries of classes) are weaker representatives and are prone to misclassification. The circle is essentially the sphere of the neighborhood of x . Typically, the confidence $\phi_i(x)$ is inversely proportional to the distance (s_i) of x from the centroid of cluster ω_i .

KNN has the advantage of not respecting fixed boundaries between classes. In the example of Figure 2.2.1, the KNN classifier would label x as '5' with a score of $\frac{3}{6} = 0.5$, as four of the samples in the K-neighborhood ($k=6$) are labeled ω_5 by the recognizer (some erroneously).

Since potentially, a test pattern can fall anywhere in a feature space, we can precompute the class that will be assigned to each point in the feature space. If the KNN rule is used, the boundaries between classes will have a complex shape. Further, the boundaries of classes will change based on the size of the neighborhood k .

Computation of such boundaries can be done ahead of time. For example applying Voronoi diagram to partition the space will require $O(k+\log n)$ time and $O(k^2n)$ space, where n is the number of prototypes and k is the number of neighbors considered. When a test pattern falls at a particular point in the feature space a table lookup can determine its class. The table can prestore the class assignments for all points in feature space for various values of k .

A typical binary feature vector of size 512 requires a feature space of size 2^{512} . A table that specifies a class for each point in the feature space will need to be 2^{512} as well. Hence, while the idea is attractive for speed reasons, the tradeoff with space storage is unacceptable.

Using Bayes rule([63]) is another way to deal with the problem, especially in the cases when a large database is available. This method is very attractive in view of recent empirical evaluations showing high accuracy([22], [40]). Unfortunately, the learning efficiency suffers when the database is incomplete.

Another alternative is to employ a neural network pattern recognition system, such as a Fuzzy ARTMAP ([8]), which synthesizes fuzzy logic and Adaptive Resonance

Theory by exploiting the formal similarity between the computation of fuzzy subsets and the dynamics of category choices, search and learning. There are, however, some limitations to this approach: (i) it is order dependent, (ii) it is severely affected by noise points, (iii) it is time consuming and (iv) it lacks an objective function and stability.

2.2.2 Our Approach

We describe a method that is based on the following two ideas which appeal to common sense [6].

1. When the correctness of a classifier on a pattern x is in question, it is best to consider the performance of the same classifier on the patterns which are “similar” to x .
2. A classifier is usually accurate when the test pattern x falls close to the center of its class in feature space and prone to error when it falls near a class boundary.

KNN addresses the first observation quite adequately by evaluating the neighbors of a test pattern x . It is expected that in any small neighborhood all samples are approximately at the same distance from any cluster center and hence are all quite “similar”. While, KNN does not oppose the second observation, we have described how it does not lend itself amenable to space storage efficiency when using a look-up table.

We propose the following two-pronged approach. For the case where a sample falls in “close” proximity of a cluster center (of a class to which it will be assigned), use the

traditional method of distances to assign the class and compute its confidence score. For the case where the distances from several cluster centers are all approximately the same, we will adopt a modification of the KNN rule. Such an approach addresses both the observations made above, and at the same time does not cause prohibitive computational issues. Since we expect only a small portion of the feature space to be near the boundaries between classes, the size of the lookup table will be reasonable.

2.2.2.1 Definitions

The main purpose of this section is to describe a new approach for dealing with the patterns from the “overlapping” area . To this end we wish to consider how one defines the terms “class center”, “close proximity” and “overlapping area”.

There are many algorithms available for clustering given a set of prototypes, and there are many different approaches for identifying the “overlapping” areas([2], [13], [34], [48]). In this section, we will choose simple and intuitive ways to describe these terms. Assuming that s_i is the maximum possible score corresponding to the class ω_i , $i = 0, \dots, 9$, the “class center” of class ω_A is defined by all prototypes y , such that $\phi_A(y) \geq \frac{3}{4}s_A$ and $\phi_B(y) \leq \frac{1}{10}s_B$ for all other classes ω_B . The “overlapping” area is defined as the part of the space of all patterns X , that contains patterns $y \in X$, such that $|\phi_A(y) - \phi_B(y)| \leq \frac{s_A + s_B}{2}$, where classes ω_A and ω_B are the top two choices for pattern y . Let us note that in some cases as in the example of Figure 2.2.2, the “overlapping” area can be narrowed even further.

2.2.2.2 Novelty

For most digit recognizers the correct class choice for an unknown pattern is almost always among the two classes with highest scores. This is to say that the top two choice correctness rate is almost 100%. Based on this observation the goal of our methodology is to overturn some of the results of the recognizer in order to promote the correct class with the second highest score to become the top choice of the recognizer. In this way the top choice correct rate will be very close to the original top two choices correct rate.

2.2.2.3 GSC

Although we have tested our method with other classifiers, we will discuss our experiments with the Gradient Structural Concavity (GSC) recognizer which is described in [17]. GSC uses 512 symbolic multi-resolutional features, which measure the image characteristics at local, intermediate, and large scales. The 192 **G**radient features measure edge curvature in a neighborhood of a pixel and provide information about stroke shape on a small scale. The 192 intermediate **S**tructural features measure short stroke types which span several pixels and give useful information about stroke trajectories. The 128 **C**oncavity features are used to detect stroke relationships at long distances which can span across the image. Features at all the 3 levels, **G**, **S**, and **C** are combined in a 512 binary feature vector and a weighted k-nearest neighbor algorithm is used for classification. If the top k prototypes x_1, \dots, x_k are chosen to be the closest to an unknown pattern x , a weighted vote for each class C is computed

as follows.

$$Weight_C = \frac{1}{k} \sum_{i=1}^k c_i d(x, x_i).$$

Here c_i is 1 if prototype x_i is labeled as class C and 0 otherwise, and $d(x, x_i) = \frac{n_{11} + n_{00}/S}{512}$, where n_{11} is the number of '1' bits that match between binary feature vectors of x and x_i , and n_{00} is the number of '0' bits that match between binary feature vectors of x and x_i and S is a parameter in the interval $[1, 5]$. The classes C are then ranked according to the value of $Weight_C$. The maximum possible score s_i , corresponding to any class is equal to 0.6.

2.2.2.4 Methodology

Let us say the top two classes are ω_A and ω_B when presented with a test pattern x . A common practice [48] is to compute $|\phi_A(x) - \phi_B(x)|$. If $|\phi_A(x) - \phi_B(x)| \leq \tau$ (where τ is a pre-determined threshold, $\tau = 0.1$ for GSC), then the confusion between the two classes is deemed to be too high to return a “confident” top choice.

In our approach we consider two classes at a time, namely classes ω_A and ω_B . A sample is considered as “easy” if $|\phi_A(x) - \phi_B(x)| \gg \tau$. Clearly, the recognizer is confident about the identity of x as either ω_A or ω_B , depending on whichever is greater of the two. A sample x in the feature space falls in the “difficult” or the “overlapping” zone when $|\phi_A(x) - \phi_B(x)| \leq \tau$.

Using the confidence functions $\phi_A : X \rightarrow [0, 0.6]$ and $\phi_B : X \rightarrow [0, 0.6]$, we define a map $\Phi_{A,B} : X_{A,B} \rightarrow \mathbb{R}^2$, such that $\Phi_{A,B}(x) = (\phi_A(x), \phi_B(x))$, where $X_{A,B}$ is the set of all patterns y , for which ω_A and ω_B are the top two classes according to the classifier. Using $\Phi_{A,B}$ any pattern from $X_{A,B}$ is mapped onto a point of the square

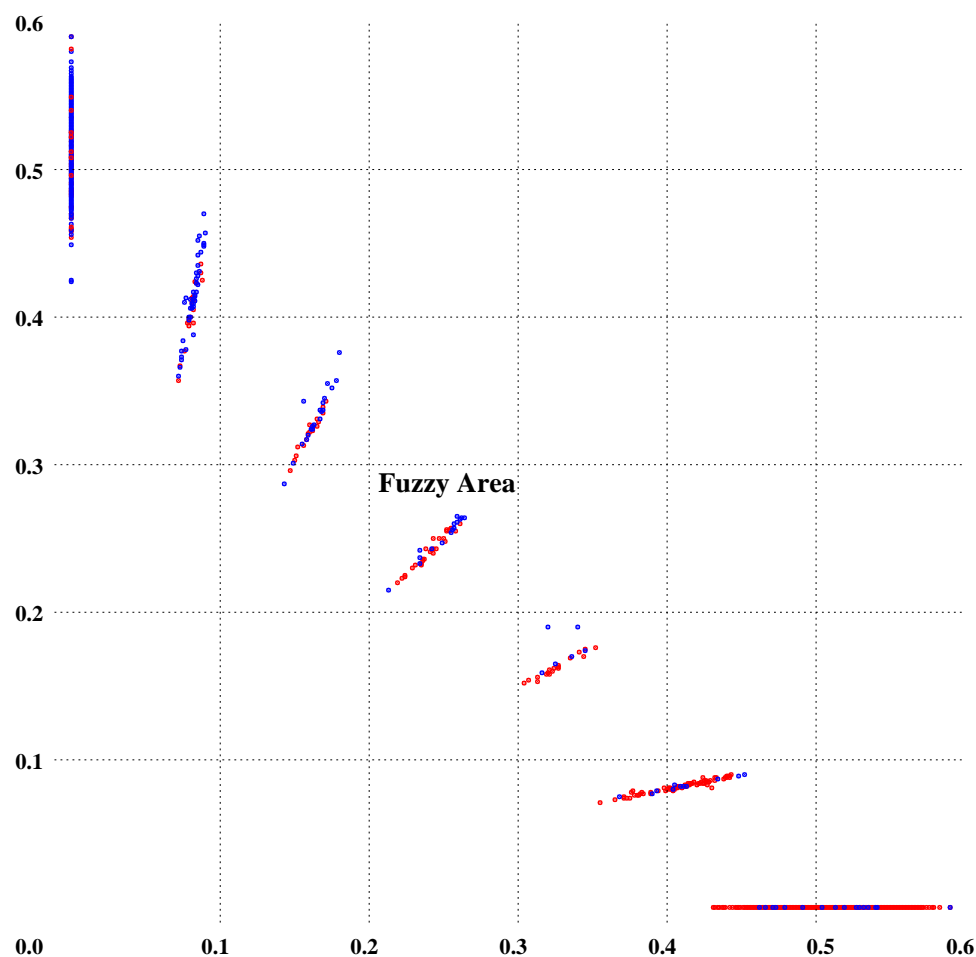


Figure 2.2.2: This is a confidence space of two classes '0' and '6'. The red dots represent the 0s and the blue dots represent the 6s. The GSC digit recognizer gives a confidence in the range of $[0,0.6]$. The truly "good" quality samples of '0' receive a score of 0.6 and the probability of misclassification is small as noted from the very few blue dots mingled with red dots at these confidence levels ($\phi_0 \approx 0.6$ and $\phi_6 \approx 0.0$). The reverse is also true. Maximum misclassifications occur when $\phi_0 \approx 0.3$ and $\phi_6 \approx 0.3$. This qualifies as the fuzzy area. These are the samples that we wish to reclassify.

$[0, 0.6]^2$, and the intersection of the “overlapping” area and $X_{A,B}$ is projected onto an area in that square, we would like to call “fuzzy”. Figure 2.2.2 shows the example of two classes ω_0 and ω_6 with τ is 0.1 ($A \equiv 0$ & $B \equiv 6$). It is to be noted that when $[\phi_A(x) \gg \phi_B(x)]$ OR $[\phi_B(x) \gg \phi_A(x)]$ then the misclassifications of such patterns are rare. It is also interesting to note how the misclassifications increase as one approaches the fuzzy area from both extremes.

Another interesting insight we gain from Figure 2.2.2 is the way the samples align in confidence space. When $\phi_6(x) \in [0.4 \dots 0.6]$ AND $\phi_0(x) \in [0.0 \dots 0.1]$ the samples in confidence space are all strongly leaning towards being a ω_6 (digit ‘6’) by aligning parallel to the x-axis. Similarly, when $\phi_0(x) \in [0.4 \dots 0.6]$ AND $\phi_6(x) \in [0.0 \dots 0.1]$ the samples in confidence space are all strongly leaning towards being a ω_0 (digit ‘0’) by aligning parallel to the y-axis. When $\phi_0(x) \in [0.2 \dots 0.3]$ AND $\phi_6(x) \in [0.2 \dots 0.3]$ the samples in confidence space represent the fuzzy area (the area of interest). Their scatter pattern aligns at about 45° .

It is also clear from the way scores are computed that a sample x cannot receive simultaneously a very high or a very low score for both the classes. This actually helps us to narrow the fuzzy area. Hence all patterns must fall in the diagonal area of the space spanning from point (0.0,0.6) to point (0.6,0.0) at an angle of 135° . Although, when we talk of cluster centers, we are referring to the confidence space, the notion of a cluster does hold meaning in the confidence space as well as can be seen in Figure 2.2.2. The cluster center for the class ω_6 would be in the vicinity of point (0.6,0.0) and the cluster center for ω_0 would in the vicinity of point (0.0,0.6).

There are 7 clusters shown in Figure 2.2.2. One of the clusters is ambiguous about its class identity. This is the fuzzy area. The 6 other clusters either favor the class ω_0 (each cluster representing a different style of writing, perhaps) or the class ω_6 . The

“goodness” or quality of the clusters improves as they move away from the fuzzy area.

Figure 2.2.3, Figure 2.2.4 and Figure 2.2.5 are other examples. Figure 2.2.3 involves classes ω_0 and ω_1 , and shows the projection of $X_{0,1}$ onto the square $[0, 0.6]^2$. Figure 2.2.4 involves classes ω_0 and ω_2 , and shows the projection of $X_{0,2}$ onto the square $[0, 0.6]^2$. Figure 2.2.5 involves classes ω_0 and ω_3 , and shows the projection of $X_{0,3}$ onto the square $[0, 0.6]^2$.

2.2.2.5 Algorithm

1. *Identify* the two classes returned as the top two choices by a digit recognizer (say, ω_A and ω_B) on pattern x .
2. *If* the quantity $|\phi_A(x) - \phi_B(x)| \gg \tau$
then accept the top choice as correct
else proceed to evaluate if the second choice should be promoted to the top by using a lookup table based on the modified KNN as described in the next section.

2.2.2.6 Analysis

Figure 2.2.6 shows samples of images that will be confidently recognized as a ‘0’ in the top choice by a traditional recognizer (such as the GSC recognizer [17]) and needs no adjustment to the classification.

Figure 2.2.7 shows samples of images that will be confidently misclassified as a ‘6’ in

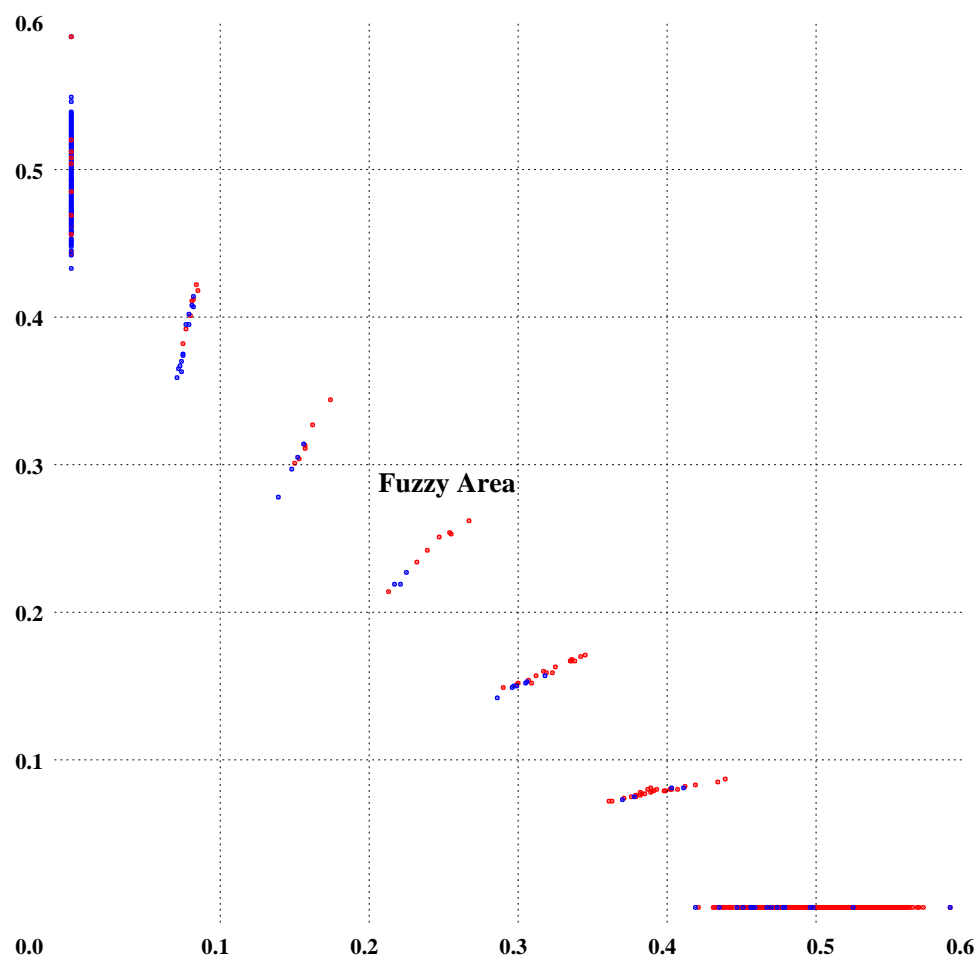


Figure 2.2.3: This is a confidence space of two classes '0' and '1'. The red dots represent the 0s and the blue dots represent the 1s. The GSC digit recognizer gives a confidence in the range of $[0, 0.6]$. The truly "good" quality samples of '0' receive a score of 0.6 and the probability of misclassification is small as noted from the very few blue dots mingled with red dots at these confidence levels ($\phi_0 \approx 0.6$ and $\phi_1 \approx 0.0$). The reverse is also true. Maximum misclassifications occur when $\phi_0 \approx 0.3$ and $\phi_1 \approx 0.3$. This qualifies as the fuzzy area. These are the samples that we wish to reclassify.

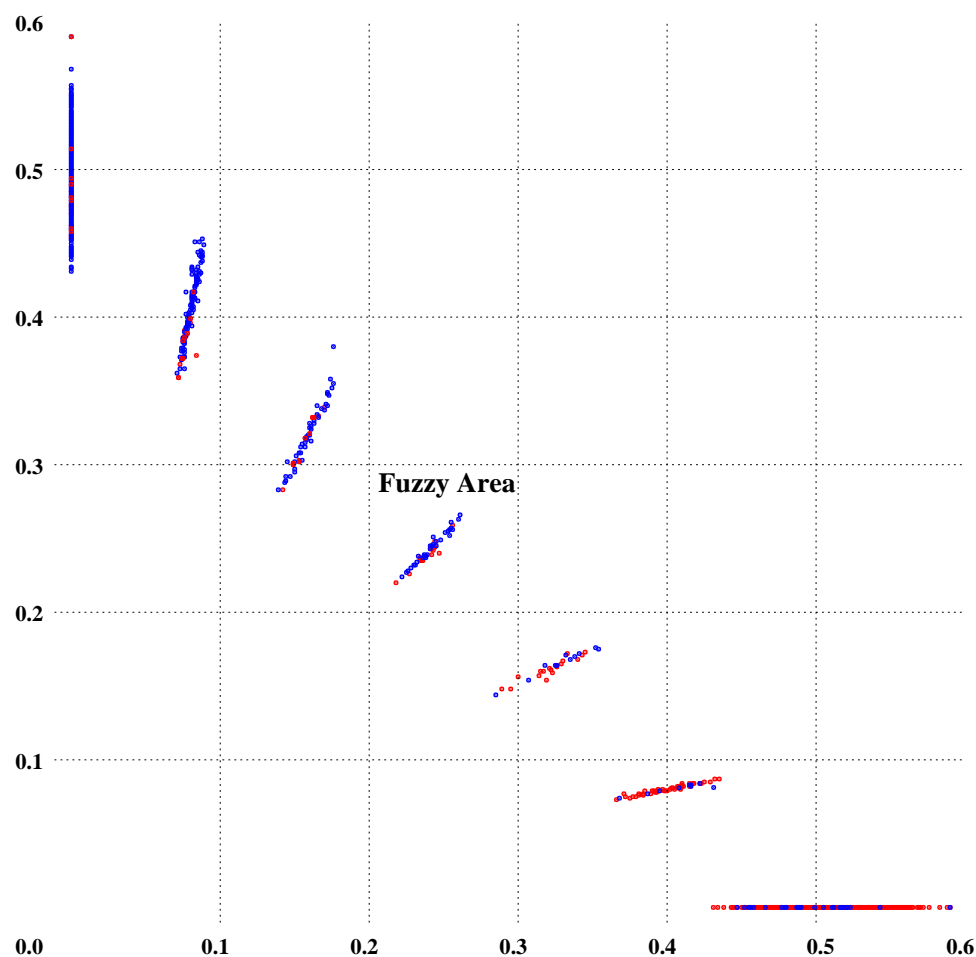


Figure 2.2.4: This is a confidence space of two classes '0' and '2'. The red dots represent the 0s and the blue dots represent the 2s. The GSC digit recognizer gives a confidence in the range of $[0,0.6]$. The truly "good" quality samples of '0' receive a score of 0.6 and the probability of misclassification is small as noted from the very few blue dots mingled with red dots at these confidence levels ($\phi_0 \approx 0.6$ and $\phi_2 \approx 0.0$). The reverse is also true. Maximum misclassifications occur when $\phi_0 \approx 0.3$ and $\phi_2 \approx 0.3$. This qualifies as the fuzzy area. These are the samples that we wish to reclassify.

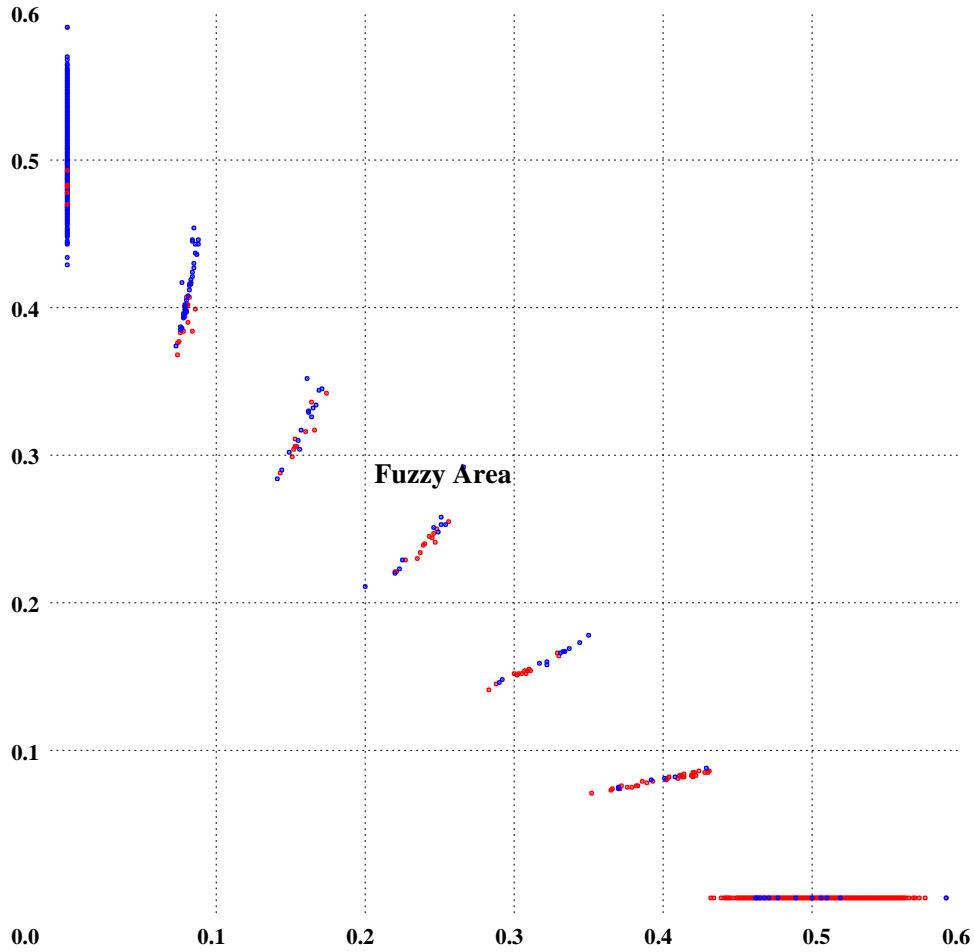


Figure 2.2.5: This is a confidence space of two classes '0' and '3'. The red dots represent the 0s and the blue dots represent the 3s. The GSC digit recognizer gives a confidence in the range of $[0, 0.6]$. The truly “good” quality samples of '0' receive a score of 0.6 and the probability of misclassification is small as noted from the very few blue dots mingled with red dots at these confidence levels ($\phi_0 \approx 0.6$ and $\phi_3 \approx 0.0$). The reverse is also true. Maximum misclassifications occur when $\phi_0 \approx 0.3$ and $\phi_3 \approx 0.3$. This qualifies as the fuzzy area. These are the samples that we wish to reclassify.

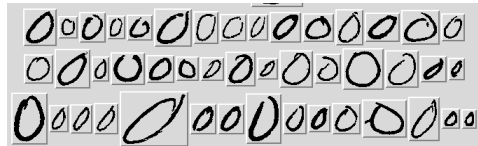


Figure 2.2.6: A sampling of “good” quality 0s. These are samples that receive $\phi_0 \approx 0.6$ and $\phi_6 \approx 0.0$ (Figure 2.2.2.)

the top choice by a traditional recognizer (such as the GSC) and cannot be readily rectified.

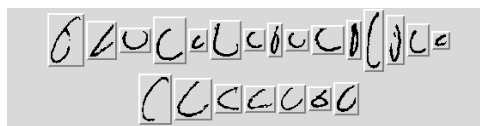


Figure 2.2.7: A sampling of “poor” quality 0s. These are samples that receive $\phi_0 \approx 0.0$ and $\phi_6 \approx 0.6$ (Figure 2.2.2.)

Figure 2.2.8 shows samples of images that will be confidently recognized as a ‘6’ in the top choice by a traditional recognizer (such as the GSC) and needs no adjustment to the classification.

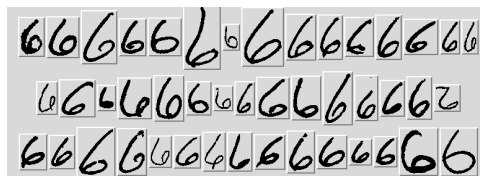


Figure 2.2.8: A sampling of “good” quality 6s. These are samples that receive $\phi_6 \approx 0.6$ and $\phi_0 \approx 0.0$ (Figure 2.2.2.)

Figure 2.2.9 shows samples of images that will be confidently misclassified as a ‘0’ in the top choice by a traditional recognizer (such as the GSC) and cannot be readily rectified.

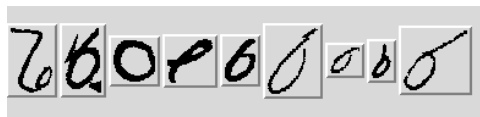


Figure 2.2.9: A sampling of “poor” quality 6s. These are samples that receive $\phi_6 \approx 0.0$ and $\phi_0 \approx 0.6$ (Figure 2.2.2.)

Figure 2.2.10 shows samples of images that will be confusing to tell apart as a ‘0’ or ‘6’ by a traditional digit recognizer such as GSC. These are the cases which we wish

to re-classify using our methodology. These are the samples that fall in the fuzzy area of the confidence space shown in Figure 2.2.2.

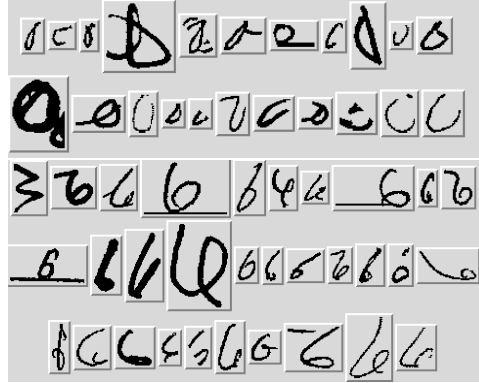


Figure 2.2.10: A sampling of “difficult” quality 0s and 6s. These are samples that receive $\phi_6 \approx 0.3$ and $\phi_0 \approx 0.3$ (Figure 2.2.2.)

Figure 2.2.11 shows the new class boundary after re-classification based on the modified KNN for the samples in the fuzzy area (Figure 2.2.2). A straight classification based on the value $\phi_i(x)$ being greater or less than a threshold (τ) would draw the class boundary as a straight line right through the middle of the fuzzy area at a 45° angle as shown. The modified KNN as will be described in the next section generates a complex boundary which does not eliminate all the classification errors, in some cases it even creates additional errors in the process of correcting some errors, nevertheless the method makes an overall improvement.

Figure 2.2.12 and Figure 2.2.13 show more examples. Figure 2.2.12 involves classes ω_0 and ω_2 and shows a new class boundary between 0s and 2s. Figure 2.2.13 involves classes ω_0 and ω_3 , and shows a new class boundary between 0s and 3s.

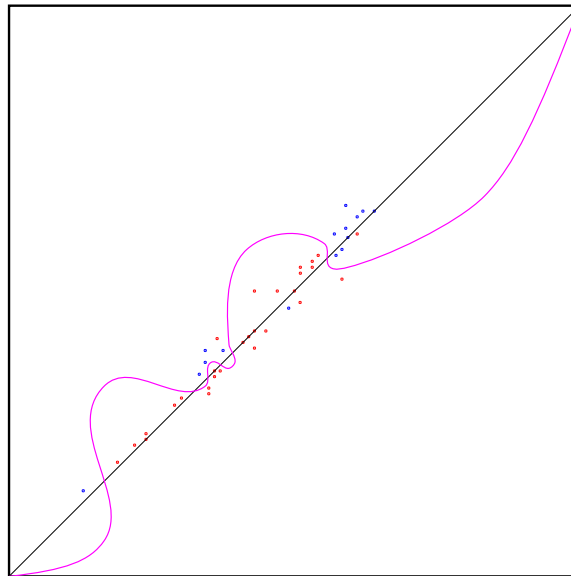


Figure 2.2.11: Re-defining the class boundary between 0s and 6s in the fuzzy area of Figure 2.2.2

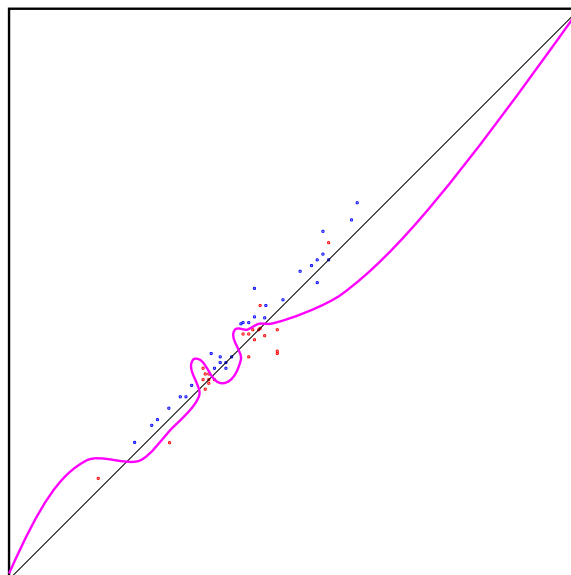


Figure 2.2.12: Re-defining the class boundary between 0s and 2s in the fuzzy area of Figure 2.2.4

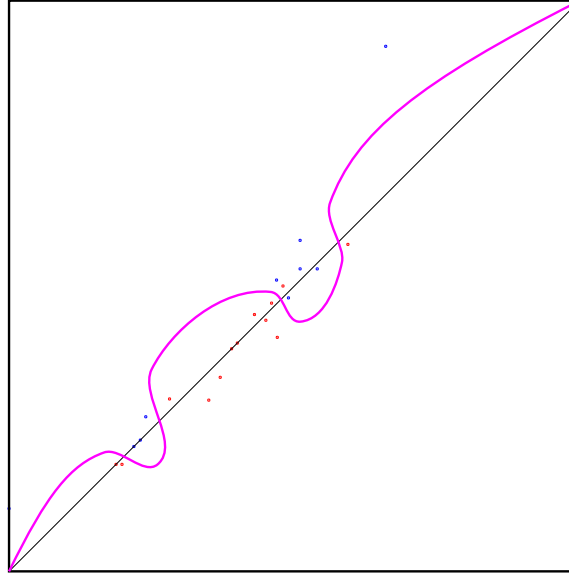


Figure 2.2.13: Re-defining the class boundary between 0s and 3s in the fuzzy area of Figure 2.2.5

2.2.3 Topology

Let us start this section by repeating the assumption that the scores for each class, returned by the recognizer, are real numbers in interval $[0, 1]$. Let U be a “blind” set of digit patterns. We will assume that all classes are equally represented in U . Let V be a subset of U that contains all patterns in U that are classified correctly by the classifier. We will use the set V to simulate the fuzzy boundaries between classes. V is partitioned into a disjoint union of subsets V_0, V_1, \dots, V_9 , where each subset V_i contains those digit patterns in V , that represent class ω_i . The score, associated with class ω_i is the highest in the confidence vector returned by the classifier, i.e

$$V_i = \{x \in V : \phi_j(x) < \phi_i(x) \text{ for all } j \neq i\}.$$

Let $f_i(x)$ be the characteristic function of the set V_i , i.e.

$$f_i(x) = \begin{cases} 1 & \text{if } x \in V_i \\ 0 & \text{otherwise.} \end{cases}$$

Let $W = [0, 1]_0 \times [0, 1]_1$, be a Cartesian product of 2 copies of the interval $[0, 1]$. Here $[0, 1]_i$ denotes the i^{th} -copy of the interval $[0, 1]$. We can convert W into a topological space by introducing one of the following two topologies: Finite or Euclidean.

1. *Finite Topology:*

Assume that each interval $[0, 1]_i$ is partitioned into finitely many disjoint subintervals $W_i^1, W_i^2, \dots, W_i^n$. Then the set \mathcal{B} of neighborhoods $W_0^{i_0} \times W_1^{i_1}$ forms a base of a topology in W . Let us enumerate the elements of \mathcal{B} , i.e $\mathcal{B} = \{B_1, B_2, \dots, B_M\}$, where $M = n^2$.

Let $W(i_1, i_2, \dots, i_j)$ be the union of the sets $B_{i_1}, B_{i_2}, \dots, B_{i_j}$, i.e

$$W(i_1, i_2, \dots, i_j) = B_{i_1} \cup B_{i_2} \cup \dots \cup B_{i_j}.$$

Let \mathcal{M} be the family of all such sets $W(i_1, i_2, \dots, i_j)$, i.e

$$\mathcal{M} = \{W(i_1, i_2, \dots, i_j) : j \in \{1, 2, \dots, M\}, 1 \leq i_1 < i_2 < \dots < i_j \leq M\}.$$

Then \mathcal{M} is a topology of W .

2. *Euclidean Topology:*

In this case we will consider W with the topology induced by the Euclidean topology of the space \mathbb{R}^2 , i.e the topology that is generated by the collection of all balls in the set W [54].

Next step is to construct a σ -algebra, that includes the topology of the set W . We consider two cases.

In the finite case, it is apparent that $B_0 \cup B_1 \cup \cdots \cup B_M = W$. Hence $W \in \mathcal{M}$. Moreover, since no two elements of the base \mathcal{B} intersect, if the set $A \in \mathcal{M}$, then $A^c \in \mathcal{M}$, where A^c denotes the complement of A in W . Also, the union of any two sets from \mathcal{M} is a set from \mathcal{M} . This, together with the fact that \mathcal{M} is finite, implies that \mathcal{M} is a σ -algebra [66].

In the Euclidean case, we define \mathcal{M} to be Borel σ -algebra, generated by the Euclidean topology.

Now, given two classes ω_i and ω_j ($i < j$) we will introduce a family of positive σ -additive measures for the σ -algebra \mathcal{M} . We define the function $\mu(\alpha, \beta) : \mathcal{M} \rightarrow \mathbb{R}$ as follows:

$$(2.2.1) \quad \mu(\alpha, \beta)(A) = \sum_{\{x \in V_i \cup V_j : \Phi_{i,j}(x) \in A\}} \alpha f_i(x) + \beta f_j(x)$$

for all $A \in \mathcal{M}$. Here $f_i(x)$ is the characteristic function of the set V_i ; $f_j(x)$ is the characteristic function of the set V_j ; and α and β are positive numbers. The values of the function $\mu(\alpha, \beta)$ are non-negative. Moreover, a function thus defined has the property, that if $A_1 \in \mathcal{M}$ and $A_2 \in \mathcal{M}$, such that $A_1 \cap A_2 = \emptyset$, then $A_1 \cup A_2 \in \mathcal{M}$ and

$$\mu_i(A_1 \cup A_2) = \mu_i(A_1) + \mu_i(A_2).$$

Indeed, if $x \in V$ is in $A_1 \cup A_2$, then $x \in A_1$ or $x \in A_2$, since $A_1 \cap A_2 = \emptyset$, implying

$$\begin{aligned} \mu(\alpha, \beta)(A_1 \cup A_2) &= \\ &= \sum_{\{x \in V_i \cup V_j : \Phi_{i,j}(x) \in A_1 \cup A_2\}} \alpha f_i(x) + \beta f_j(x) = \\ &= \sum_{\{x \in V_i \cup V_j : \Phi_{i,j}(x) \in A_1\} \cup \{x \in V_i \cup V_j : \Phi_{i,j}(x) \in A_2\}} \alpha f_i(x) + \beta f_j(x) = \\ &= \sum_{\{x \in V_i \cup V_j : \Phi_{i,j}(x) \in A_1\}} \alpha f_i(x) + \beta f_j(x) + \sum_{\{x \in V_i \cup V_j : \Phi_{i,j}(x) \in A_2\}} \alpha f_i(x) + \beta f_j(x) = \end{aligned}$$

$$= \mu(\alpha, \beta)(A_1) + \mu(\alpha, \beta)(A_2).$$

So, we can conclude that $\mu(\alpha, \beta)$ is a family of positive σ -additive measures [66]. Let us give special names of two members of this family. We introduce

$$\mu_i = \mu(1, 0) \quad \text{and} \quad \mu_j = \mu(0, 1).$$

2.2.4 Fuzzy Sets

In the finite case, we have a set of neighborhoods

$$\mathcal{B} = \{B_1, B_2, \dots, B_M\}.$$

For every $y \in W$, there exists a unique neighborhood B_y , containing y . In the Euclidean case, given a function $R : W \rightarrow \mathbb{R}$, the set of neighborhoods of y includes all balls of the forms

$$B_y = \{z \in W : \text{dist}(y, z) \leq R(y)\}.$$

For every $y \in W$, let $\chi_y : W \rightarrow \mathbb{R}$ be the characteristic function of its neighborhood B_y . Let us also define the family of functions

$$\begin{aligned} \varphi(\alpha, \beta)(y) &= \int_W \chi_y(z) d\mu(\alpha, \beta)(z) = \\ &= \alpha |\{x \in V_i : \Phi_{i,j}(x) \in B_y\}| + \beta |\{x \in V_j : \Phi_{i,j}(x) \in B_y\}|, \end{aligned}$$

as well as the functions

$$\varphi_i(y) = \int_W \chi_y(z) d\mu_i(z) = |\{x \in V_i : \Phi_{i,j}(x) \in B_y\}|$$

and

$$\varphi_j(y) = \int_W \chi_y(z) d\mu_j(z) = |\{x \in V_j : \Phi_{i,j}(x) \in B_y\}|.$$

Let us observe, that all functions $\varphi(a, \beta)(y)$, and in particular $\varphi_i(y)$ and $\varphi_j(y)$, are nonnegative. Given a positive number α , let

$$\eta_i(y) = \begin{cases} \frac{\varphi_i(y)}{\varphi(1, \alpha)(y)} & \text{if } \varphi(1, \alpha)(y) \neq 0 \\ 0.5 & \text{otherwise} \end{cases}$$

and

$$\eta_j(y) = \begin{cases} \frac{\varphi_j(y)}{\varphi(\frac{1}{\alpha}, 1)(y)} & \text{if } \varphi(\frac{1}{\alpha}, 1)(y) \neq 0 \\ 0.5 & \text{otherwise} \end{cases}.$$

Let us observe, that for every $y \in W$, $0 \leq \eta_i(y) \leq 1$, $0 \leq \eta_j(y) \leq 1$ and $\eta_i(y) + \eta_j(y) = 1$.

Now, we are in a position to define the fuzzy sets corresponding to the different classes. Let x be an unknown pattern, and classes ω_i and ω_j correspond to the two highest entries in the confidence vector

$$\Phi(x) = (\phi_0(x), \phi_1(x), \dots, \phi_9(x)).$$

Let τ be the threshold that determines the overlapping area. We can define the fuzzy membership functions as follows:

$$\begin{aligned} \xi_i(x) &= \begin{cases} \eta_i(\Phi_{i,j}(x)) & \text{if } x \text{ is in the overlapping area} \\ \phi_i(x) & \text{otherwise} \end{cases} \\ \xi_j(x) &= \begin{cases} \eta_j(\Phi_{i,j}(x)) & \text{if } x \text{ is in the overlapping area} \\ \phi_j(x) & \text{otherwise} \end{cases} \\ \xi_i(x) &= \begin{cases} 0 & \text{if } x \text{ is in the overlapping area} \\ \phi_i(x) & \text{otherwise} \end{cases}. \end{aligned}$$

Each of these functions $\xi_i(x)$ is the membership function of the corresponding class ω_i . We say that the unknown pattern x belongs to the class ω_i , if $\xi_j(x) < \xi_i(x)$ for

all $j \neq i$. Having defined our membership functions $\xi_i(x)$, we are able to partition the space W into class regions and boundaries, that separates them in the following way. The set

$$\Omega_i = \{x : \xi_i(x) > \xi_j(x) \text{ for all } j \neq i\},$$

$i = 0, 1, \dots, 9$ is the region, all points of which are in class ω_i . The set

$$D_{ij} = \{x : \xi_i(x) = \xi_j(x) \text{ and } \xi_i(x) > \xi_k(x) \text{ for all } k \neq i, j\},$$

is the implicitly defined surface, that is a part of the boundary between classes ω_i and ω_j .

An illustration of these concepts (when only two classes are involved) is provided in Figures 2.2.2, 2.2.3, 2.2.4, and 2.2.5. Initially, the recognizer uses a linear discriminant, namely the line $y = x$, in the projection space. Figures 2.2.11, 2.2.12, etc. show the new nonlinear discriminant between the classes using the algorithm as described.

2.2.5 Bayesian View of the Membership Functions

Let x be a pattern in the overlapping area and ω_A and ω_B are the two classes with the highest and second highest scores in the confidence vector. Let s_A and s_B be the scores corresponding to ω_A and ω_B . Let us compute the *a posteriori* probabilities for the classes ω_A and ω_B being the true class for pattern x , given scores s_A and s_B of the two top class choices ω_A and ω_B .

$$P(\omega_A/s_A, s_B) = \frac{P(s_A, s_B/\omega_A)P(\omega_A)}{P(s_A, s_B/\omega_A)P(\omega_A) + P(s_A, s_B/\omega_B)P(\omega_B)}.$$

Without any prior information, we can assume that both classes have equal chance of being the truth, so we can reasonably assume that $P(\omega_A) = P(\omega_B)$. Then

$$\begin{aligned}
 P(\omega_A/s_A, s_B) &= \frac{P(s_A, s_B/\omega_A)}{P(s_A, s_B/\omega_A) + P(s_A, s_B/\omega_B)} = \\
 &= \frac{P(s_A/\omega_A)P(s_A/s_B, \omega_A)}{P(s_A/\omega_A)P(s_A/s_B, \omega_A) + P(s_B/\omega_B)P(s_B/s_A, \omega_B)} = \\
 &= \frac{P(s_A/\omega_A)}{P(s_A/\omega_A) + \alpha P(s_B/\omega_B)},
 \end{aligned}$$

where $\frac{P(s_B/s_A, \omega_B)}{P(s_A/s_B, \omega_A)}$. In the same way we compute

$$P(\omega_B/s_A, s_B) = \frac{P(s_B/\omega_B)}{\frac{1}{\alpha}P(s_A/\omega_A) + P(s_B/\omega_B)}.$$

By definition, the probability records the frequency of the results of an event over many trials. However, given a pattern x , the classifier returns the same confidence vector over and over, no matter how many times the recognition process is repeated. In order to deal with this problem, i.e to compute $P(s_A/\omega_A)$, we will approximate the probability by using patterns y , which are “close” to pattern x in the “eyes” of the recognizer, i.e $\Phi_{A,B}(y)$ falls in a given neighborhood B_x of $\Phi_{A,B}(x)$. Then,

$$P(s_A/\omega_A) = \frac{\sum_{y \in V_A \cup V_B : \Phi_{A,B}(y) \in B_x} f_A(y)}{|V_A \cup V_B|} = \frac{\mu_A(B_X)}{|V_A \cup V_B|}.$$

Similarly,

$$P(s_B/\omega_B) = \frac{\mu_B(B_X)}{|V_A \cup V_B|}$$

and therefore

$$P(\omega_A/s_A, s_B) = \frac{\mu_A(B_X)}{\mu_A(B_X) + \alpha\mu_B(B_X)} = \xi_A(x)$$

and

$$P(\omega_B/s_A, s_B) = \frac{\mu_B(B_X)}{\mu_B(B_X) + \frac{1}{\alpha}\mu_A(B_X)} = \xi_B(x).$$

2.3 Experiments

We have already described how the fuzzy boundaries change results only for those samples that fall in the proximity of class boundaries. Hence if the the original recognizer has a high recognition accuracy, our method can only marginally affect the results. GSC digit recognizer is such a recognizer with very high accuracy. The method described in this section successfully reduced its error rate from 2.85% to 1.96%, i.e by 0.89%, which is more than 30% from the initial error of 2.85% (Table 2.3.1). To illustrate the significance of the improvement let us look at a simple application. Assuming we use GSC to recognize 5 digits in handwritten ZIP codes, taken from the US mail stream. GSC recognizes correctly the entire ZIP code in 86.54% of the cases and the method, described in this section will recognize correctly the entire ZIP code in 90.58% of the cases.

The training set U (50000 images) and the testing set (30398 images) were created using digit samples extracted from the US mail stream. There are two reasons why we use our own database. First, GSC achieves almost 100% correct rate on databases available publicly such as NIST and CENPARMI. Second, all classes are equally represented in the training set U which is not the case with other databases.

In the experiment, we use $\tau = 0.1$ to determine the overlapping area. Then the confidence space is partitioned as described in the previous sections into squares using $n = 100$. Since the training is off-line, we precompute the measures of any of the squares in the fuzzy area and then store them in a lookup table. Once the recognizer returns the confidence vector $\Phi(x)$ of an unknown pattern, a check is performed to see if the pattern is in the “overlapping area”. If it is in the overlapping area, the corresponding entries in the lookup table are used to determine the values of the fuzzy membership functions for the pattern and its classification is performed.

One way to evaluate the performance of two recognizers is to compare the corresponding percentage correctly classified samples (Tables 2.3.1, 2.3.2).

test set	GSC		Fuzzy	
	correct	%	correct	%
30398	29532	97.15	29807	98.04

Table 2.3.1: Improvements to GSC recognizer using Fuzzy boundaries

Actually, the best way to evaluate the performance of two classifiers is to compare their graphs on a REJECT vs ERROR scale (Figure 2.3.1).

Figure 2.3.1 shows the graph of GSC recognizer and the fuzzy method on a REJECT vs ERROR scale in order to compare their performances. Table 2.3.3, shows some selected points from the graphs of both recognizers (Figure 2.3.1) in order to compare the correct rates of classifiers given the reject rate.

class	test set	GSC		Fuzzy	
		correct	%	correct	%
0	6986	6833	97.81	6852	98.08
1	6554	6496	99.12	6496	99.12
2	3954	3759	95.07	3845	97.24
3	2238	2170	96.96	2186	97.68
4	3191	3070	96.21	3123	97.87
5	1556	1497	96.21	1513	97.24
6	1474	1442	97.83	1449	98.30
7	1445	1391	96.26	1412	97.72
8	2049	1945	94.93	2001	97.61
9	951	929	97.69	930	97.93

Table 2.3.2: Improvements to GSC recognizer using Fuzzy approach for any of the classes

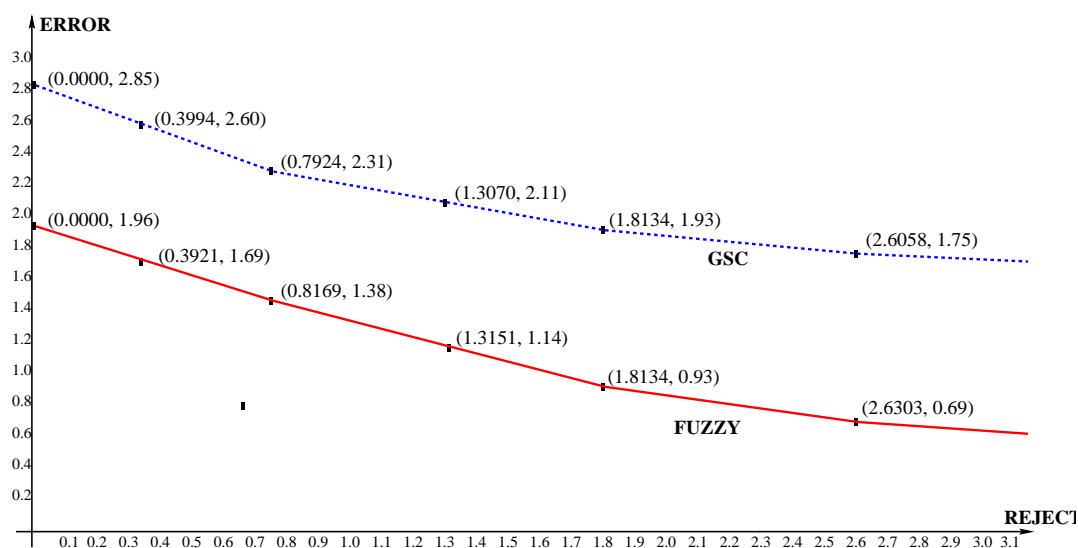


Figure 2.3.1: This figure shows how the performances of GSC recognizer and the proposed fuzzy method are compared on a REJECT vs ERROR scale. Based on the fact the graph of the first fuzzy method is upon the graph of GSC, we can conclude that the overall performance of our method is better the one of GSC recognizer

GSC		Fuzzy	
reject %	error %	reject %	error %
0.0000	2.85	0.0000	1.96
0.3994	2.60	0.3921	1.69
0.7924	2.31	0.8169	1.38
1.3070	2.11	1.3151	1.14
1.8134	1.93	1.8134	0.93
2.6303	1.75	2.6058	0.69

Table 2.3.3: This table shows some selected points from the graphs of both recognizers on a REJECT vs ERROR scale in order to compare the correct rates of classifiers given a reject rate and shows improvement of the fuzzy method over GSC at those reject rates

At the end, we have performed the similar experiment with other available recognizer and obtain similar results (Table 2.3.4).

test set 30398	Polynomial		Fuzzy	
	correct	%	correct	%
	27744	91.27	28155	92.62
	Gradient		Fuzzy	
	correct	%	correct	%
	28671	94.32	29039	95.53
	Curvature		Fuzzy	
	correct	%	correct	%
	29191	96.03	29513	97.09

Table 2.3.4: Improvements to others recognizer using fuzzy method

2.4 Summary

First, let us discuss the implementation issues. In the case of finite topology, a lookup table of size $O(n^2)$ is created off-line in $O(n|U|)$ time, where $|U|$ is the size of the training set and n is the number used to partition the confidence space. The fuzzy method transforms the confidence vector into a new vector, containing the values of

the fuzzy membership function in $O(1)$ time. In the Euclidean case, we store the confidence vectors of all samples in the training set U ($O(|U|)$ space is required) and then the needed measures are computed on-line in at most $O(|U|)$ time. So, this version of the fuzzy method is slower, but it is more accurate.

Second, let discuss a modification of the fuzzy algorithm. Let us consider 3 or more classes: x is a unknown pattern, ω_A is the class with highest score in the confidence vector, and classes satisfy the “overlapping” condition, i.e $|\varphi_A(x) - \varphi_B(x)| < \tau$, $|\varphi_A(x) - \varphi_C(x)| < \tau$. In other words, the pattern is around the boundaries of three or more classes. Situations like this occur seldom, but still require some attention. A modification of the “fuzzy” method described can easily deal with situations like this. Using the confidence vector, we project space X of all patterns in a cube in \mathbb{R}^{10} , then we introduce a topology similar to the one(finite or Euclidean) described before and then use the following family of measures to define the fuzzy membership functions:

$$\mu(\alpha_0, \alpha_1, \dots, \alpha_9)(A) = \sum_{\{x \in V: \Phi(x) \in A\}} \sum_{i=0}^9 \alpha_i f_i(x),$$

under the condition that $\alpha_i = 0$ if $|\varphi_A(x) - \varphi_B(x)| \geq \tau$.

Another issue of importance is the way in which the “overlapping” area is determined. In this section we use one very simple and intuitive way to describe the “overlapping” area. In the case of GSC, this was very effective (Figure 2.2.2). In cases where other classifiers are used, this issue is more complex and requires more investigation. There some interesting problems have arisen: should more than one threshold be used to determine the overlapping area; what is the relation between the threshold and the size and density of the overlapping area; how does one choose the optimal neighborhood for each point in the fuzzy area and what topology corresponds to this set of neighborhoods; what is the relation between the misclassified patterns that have been classified correctly by our method and correctly classified patterns misclassified

by our method. These and other open issues we plan to investigate in the future.

Finally, the method described is able to improve a recognition system while treating it as a black box. Further, the only resource required to enable this method is a large database of samples for retraining. More samples lead to the better results.

Chapter 3

Applications-Classifier Combination

For complex pattern recognition tasks, designing a good classifier is a tough problem. In many practical applications combination of simple classifiers can be an attractive alternative. In many situations there is more than one classifier available, so instead of ignoring them we can combine their results.

The idea of combining the decisions of multiple experts has been explored by many researchers [38], [43], [72], [3], [1], [28], [83]. For experts of Type I for which only labels are available, voting algorithms are used ([38], [72], [1]). Label rankings are used with experts of Type II ([1], [3], [28]). In case of Type III experts with measurement values interpreted as posteriori probabilities, a Bayesian technique is often applied for experts assemblies ([83], [60]). If the expert output is interpreted as fuzzy membership or evidence values, fuzzy rules ([9],[10]) and Dempster-Shafer approach ([83],[65], [53]) are used. Also there are cases of expert combination, where the output of the expert

is used as a new feature and a new expert (neural network) is built to perform the combination ([81], [31], [39]).

In this chapter, we focus on a simplified version of Type III multi-expert combination. Each expert uses its own representation, i.e. measurements extracted from the input pattern are unique to each expert. We create a new accuracy measurement scale, uniform for all experts and use it as a basis for the expert combination. We develop an optimal combination scheme which requires extremely large amount of training data as well as memory space. So we introduce an empirical scheme to approximate the optimal scheme so that the requirements on the amount of training data and memory is practical. We first develop the accuracy measure which is common for all experts. Next, we characterize each expert with a family of accuracy maps. Next, we build accuracy combination maps with a special synthetic function. Finally, we construct synthetic accuracy maps for the combined confidences. We also propose a simplified combination scheme which requires less training data while sacrificing some accuracy. We finally discuss the optimal rejection threshold for the final recognition decision.

3.1 Confidence Combination Methods in Multi-Expert Systems

Many classifiers are able to supply confidence information from the measurement level. Bayes' classifier supplies the a posteriori probabilities as confidence measurements. Various distance related classifiers use the distance between a test pattern and template/prototype patterns as confidence measurement.

In this section, we assume that the given experts only supply the top choice and the corresponding confidence. Specifically, let $E^{(n)}$ represent expert (classifier) n , where $n = 1, 2, \dots, N$, and N is the total number of experts. $\Lambda = \{C^{(1)}, C^{(2)}, \dots, C^{(M)}\}$, are mutually exclusive and exhaustive set of class labels. $E^{(n)}(\mathbf{x})(E_1^{(n)}(\mathbf{x}), E_2^{(n)}(\mathbf{x}))$ means that expert n assigns the unknown pattern \mathbf{x} to class label $E_1^{(n)}(\mathbf{x}) \in \Lambda$ with confidence $E_2^{(n)}(\mathbf{x}) \in [0, 1]$.

Our task is to approximate the conditional accuracy distribution function for a multi-expert system $P(E_1^{(i)} | (E_1^{(1)}, E_2^{(1)}), (E_1^{(2)}, E_2^{(2)}), \dots, (E_1^{(N)}, E_2^{(N)}))$, $i = 1, 2, \dots, N$ by using the given training samples.

3.1.1 Accuracy Approximation

Recognizers usually differ in their confidence measures. A judicious combination of these measures can be made only when they are on the same scale. We will develop a way to transform confidence measures from different recognizers to a “special” confidence, we call accuracy. For a given set of “large enough” as well as “properly representative” training samples, the accuracy measure can be obtained by transforming the confidence using sufficient training samples. Usually, the accuracy is an increasing function of confidence.

Let us assume that the training set is adequately representative. Now let us discuss the accuracy approximation by utilizing the training samples. We assume that the confidence measures of all classifiers are continuous, that is, the confidence values can be any point in $[0, 1]$. There are many methods that can implement this transform. Here we introduce a simple and efficient method.

Let L patterns be classified to a certain class label by an expert and have confidence between $[a, b)$. Let t out of L patterns be correctly classified, then we can assign the approximate accuracy over $[a, b)$ as $\tilde{m}(r) = t/L$, $\forall r \in [a, b)$. For a given error bound ε , we claim that L has statistical sense if the probability that the approximate accuracy $\tilde{m}(r)$, $r \in [a, b)$ is within the error ε from any true accuracy value $\mu(r)$, $\forall r \in [a, b)$ is greater than $1 - \varepsilon$, that is,

$$(3.1.1) \quad Pr(\max_{r \in [a, b)} \{|\tilde{m}(r) - \mu(r)| < \varepsilon\}) > 1 - \varepsilon.$$

According to (3.1.1), the difference of two accuracy values between two adjacent representative intervals $[a, b)$ and $[b, c)$ satisfies

$$(3.1.2) \quad Pr(|\tilde{m}(r_2) - \tilde{m}(r_1)| < 2\varepsilon) > 1 - 2\varepsilon, \quad \forall r_1 \in [a, b), \forall r_2 \in [b, c).$$

At the same time, the accuracy value is within $[0, 1]$. So the number of representative intervals over $[0, 1]$ must be greater than $1/(2\varepsilon)$. Therefore, the number of training samples to approximate accuracy is of the order of $O(L/(2\varepsilon))$ for a class label. Since there are M labels, totally $O(ML/(2\varepsilon))$ samples are required to capture the accuracy characteristic of an expert with error tolerance within ε .

In the real implementation, for a given ε we are not able to estimate L according to (3.1.1). Instead, we estimate L by making the difference of two accuracy values between two adjacent representative intervals $[a, b)$ and $[b, c)$ less than 2ε , that is,

$$(3.1.3) \quad 0 \leq \tilde{m}(r_2) - \tilde{m}(r_1) < 2\varepsilon, \quad \forall r_1 \in [a, b), \forall r_2 \in [b, c).$$

3.1.2 Optimal Combination Rule

Let us assume we are given sufficient number of representative training samples so that we are able to generate the conditional accuracy distribution function for a multi-expert system, $P(E_1^{(i)} | (E_1^{(1)}, E_2^{(1)}), (E_1^{(2)}, E_2^{(2)}), \dots, (E_1^{(N)}, E_2^{(N)}))$, $i = 1, 2, \dots, N$. This “Behavior-Knowledge Space” scheme produces the optimal combination performance [31].

However, such approach needs large number of samples. Let us find the lower bound on the number of samples for the above accuracy approximation with error bound ε . There are $O(M/(2\varepsilon))$ choices for each pair $(E_1^{(i)}, E_2^{(i)})$, $i = 1, 2, \dots, N$ and thus $[M/(2\varepsilon)]^N$ choices for $\{(E_1^{(1)}, E_2^{(1)}), (E_1^{(2)}, E_2^{(2)}), \dots, (E_1^{(N)}, E_2^{(N)})\}$. For each choice, at least L samples are required. Therefore, a total number of $O(L[M/(2\varepsilon)]^N)$ samples are necessary to build the above “Behavior-Knowledge Space”. Also, the memory of $O(N[M/(2\varepsilon)]^N)$ are required to build the look-up table of the joint accuracy distribution.

3.1.3 Empirical Combination Rules

3.1.3.1 Accuracy combination functions

The accuracy combination function is of the form $F(a_1, a_2, \dots, a_n)$, where $a_i \in [0, 1]$, $1 \leq i \leq n$ are n accuracy variables. $F(\cdot)$ is supposed to be symmetric and in $[0, 1]$. Moreover, $F(\cdot)$ must satisfy the following two special properties.

$$(3.1.4) \quad F(a_1, a_2, \dots, a_n) = 1, \text{ if } a_k = 1, \exists k;$$

$$(3.1.5) \quad F(a_1, a_2, \dots, a_n) = 0, \text{ if } a_k = 0, \exists k.$$

The justification of (3.1.4) and (3.1.5) is obvious. In fact, it is never the case that $a_i = 1$ at the same time as $a_j = 0$ (theoretically). Therefore, in this case the function $F(\cdot)$ is supposed to be non-existent.

A family of functions satisfying the required conditions are:

$$(3.1.6) \quad F(a_1, a_2, \dots, a_n) = [h^{-1}(\frac{1}{n} \sum_{i=1}^n h(a_i^\lambda))]^{1/\lambda}, \quad \lambda > 0.$$

where λ is a parameter and $h(r)$, $r \in [0, 1]$ is a strictly ascending function satisfying

$$(3.1.7) \quad h(0) = -\infty;$$

$$(3.1.8) \quad h(1) = \infty.$$

Here, we consider ∞ as an existent number. Four simple examples of $h(\cdot)$ are listed as follows:

$$(3.1.9) \quad h(r) = \tan(\pi(r - 1/2)), \quad r \in [0, 1];$$

$$(3.1.10) \quad h(r) = (1/2 - |r - 1/2|)^{-1}(r - 1/2), \quad r \in [0, 1];$$

$$(3.1.11) \quad h(r) = (r - r^2)^{-1}(r - 1/2), \quad r \in [0, 1];$$

$$(3.1.12) \quad h(r) = (r - r^2)^{-1/2}(r - 1/2), \quad r \in [0, 1].$$

3.1.3.2 Combination scheme

For each classification distribution of the multi-expert system (C_1, C_2, \dots, C_N) , we build an accuracy map, called characteristic accuracy map for each expert. We note that there are M^N possibilities. So we build M^N characteristic accuracy maps in total for each expert. In practice, we might not have all the M^N possibilities, since some of permutations do not exist. The required number of training samples with error bound ε is $O(M^N L / (2\varepsilon))$. We note the required training samples is just linear

in the number required to approximate an accuracy map, thus greatly reducing the required number of training samples.

For a given distribution (C_1, C_2, \dots, C_N) , we denote a set of training samples by χ such that

(3.1.13)

$$(E_1^{(1)}(\mathbf{x}), E_1^{(2)}(\mathbf{x}), \dots, E_1^{(N)}(\mathbf{x})) = (C_1, C_2, \dots, C_N), \quad \forall \mathbf{x} \in \chi.$$

Let Γ denote the expert indices set such that $C_i = C, \forall i \in \Gamma$ and $|\Gamma| \geq 2$. Let $\tilde{E}_2^{(i)}, i \in \Gamma$ denote the characteristic accuracy maps constructed from χ and $E_1^{(i)}(\mathbf{x}) = C, \forall i \in \Gamma, \forall \mathbf{x} \in \chi$. We define the synthetic accuracy as

(3.1.14)

$$E_2^{(\Gamma)}(\mathbf{x}) = F(\tilde{E}_2^{(i)}(\mathbf{x}) : i \in \Gamma).$$

When the combination scheme is given, we can easily get the combination accuracy $\tilde{E}_2^{(\Gamma)}$ from the synthetical confidence $E_2^{(\Gamma)}$ using the data set χ . The generated accuracy map is supposed to be an ascending function of the combination confidence.

Now let us discuss the maximum memory required for the look-up table for all accuracy maps. For simplicity, we just consider the number of accuracy maps. The number of synthetic accuracy maps which combines exactly n accuracies are $M \binom{N}{n} (M-1)^{N-n}$. The total number of synthetic accuracy maps is $\sum_{n=2}^N M \binom{N}{n} (M-1)^{N-n} = M(M^N - N(M-1)^{N-1} - (M-1)^N)$. The maximum number of accuracy maps, including characteristic and synthetic maps, are $M(M^N - N(M-1)^{N-1} - (M-1)^N) + NM^N$.

3.1.3.3 Simplified combination scheme

In the original scheme, we need to build $O(M^N)$ accuracy maps for each expert. Usually, M is a large number, e.g., $M = 10$ in numerical recognition; $M = 10 + 26 \times 2 = 62$ in character recognition. So M^N can be a very large number even for $N = 2$. Thus this method still needs quite a large amount of training data. So instead of collecting training data for each specific distribution (C_1, C_2, \dots, C_N) in the original scheme, we can collect training data for each case such that $C_i = C, \forall i \in \Gamma$, where $C \in \Lambda$, and Γ is a set of expert indices such that $|\Gamma| \geq 2$. This method requires smaller set of training data, however this is gained by sacrificing the performance. This method performance is inferior to the original scheme when given sufficient representative training data.

The synthetic accuracy map construction procedure is same as the original scheme. For a special expert n , expert index set Γ including index n has 2^{N-1} choices and class label C has M choices. Thus we need to build $2^{N-1}M$ characteristic accuracy maps for each expert. Therefore $O(2^{N-1}ML/(2\varepsilon))$ limit samples are required. In return, we curtail the required number of training samples by $O((M/2)^{N-1})$ times. When $M = 2$, both schemes become identical.

Now let us discuss the maximum number of accuracy maps. The number of synthetic accuracy maps which combines exactly n accuracies are $M \binom{N}{n}$. The total number of synthetic accuracy maps is $\sum_{n=2}^N M \binom{N}{n} = M(2^N - 1 - N)$. The maximum number of accuracy maps, including characteristic and synthetic maps, are $M(2^N - 1 - N + N2^{N-1})$. In comparison with the original combination scheme, we also reduce the necessary memory by $O((M/2)^{N-1})$ times.

3.1.4 Optimal Rejection Threshold

In the final stage of making the recognition decision, we have to make one of two decisions: acceptance or rejection. There is a cost associated with both error as well as rejection. Trade-offs between the rejection and error ratio must be made. We follow the optimization objective as in [43]

$$(3.1.15) \quad R_{obj} = \min\{R_{err} + \alpha R_{rej}\}.$$

where R_{err} and R_{rej} are error ratio and rejection ratio, respectively, and $0 \leq \alpha \leq 1$ is a deterministic parameter.

A natural way to determine the recognition class is to choose the class which has the maximum accuracy upon the proposed scheme. However, we need to decide to discard or accept the recognition class according to the recognition accuracy.

Let us determine the accuracy threshold which minimizes the objective defined in (3.1.15). In fact, we can explicitly express both items R_{err} and R_{rej} as functions of a threshold θ as follows:

$$(3.1.16) \quad R_{err} = \int_{r=\theta}^1 (1-r)dr = 1/2 - \theta + \theta^2/2.$$

$$(3.1.17) \quad R_{rej} = \int_0^\theta 1dr = \theta.$$

Hence, we have

$$(3.1.18) \quad R_{obj} = \min_{\theta \in [0,1]} \{1/2 - \theta + \theta^2/2 + \alpha\theta\}$$

Taking the derivative of both sides with respect to θ , we obtain

$$(3.1.19) \quad 0 = -1 + \theta + \alpha$$

Thus the optimal accuracy threshold θ is given by:

$$(3.1.20) \quad \theta = 1 - \alpha.$$

That is, when recognition accuracy is less than the decision threshold $\theta = 1 - \alpha$ we reject the result, otherwise we accept it.

When $\alpha = 1$, rejecting an unknown pattern is equivalent to misclassification, as we would like to accept all recognition. This verifies the optimal threshold $\theta = 0$. When $\alpha = 0$, the objective is to minimize the error rate alone, so we would accept an unknown pattern only if it has recognition accuracy 1, that is, no error is made (theoretically). This again verifies the optimal threshold of $\theta = 1$.

3.2 Experimental Results

The training set used in the construction of the accuracy approximation and the testing set were created using digit samples extracted from the US mail stream. There are two reasons why we use our own database. First, the recognizers used in the combinations schemes achieve almost 100% correct rate on databases available publicly, such as NIST. Second, all classes are equally represented in the training set which is not the case with other databases. The training set contains 120,000 digit samples, and the testing set contains 30,000 digit samples.

We used the simplified method exactly as it is described in section V.3 and the thresholds given in the previous section. Table 1 shows the performance of the “polynomial” expert and the performance of the “gradient” expert. Table 2 shows the results of the combination of these two experts. As we can see, for all values of α , the combination method got significant improvement of the objective function - 23.93% for $\alpha = 0.2$;

α	polynomial			gradient		
	err(%)	rej(%)	opt	err(%)	rej(%)	opt
1/5	4.49	31.61	54.04	1.92	12.15	21.73
1/10	0.19	82.64	84.52	1.41	17.30	31.35
1/15	0.19	82.64	85.46	1.03	22.27	37.77
1/20	0.19	82.64	86.46	0.82	26.37	42.86

Table 3.2.1: Performances of polynomial and gradient experts

α	err(%)	rej(%)	opt
1/5	1.44	9.34	16.54
1/10	0.97	12.88	22.63
1/15	0.76	15.42	26.88
1/20	0.69	17.20	30.92

Table 3.2.2: Performance of polynomial-gradient combination

27.75% for $\alpha = 0.1$; 28.86% for $\alpha = 0.067$; and 27.85% for $\alpha = 0.05$ compared to the values of the objective function for the better expert. Table 3 describes the performance of the “cp” expert and the performance of “gsc” expert. Table 4 shows the results of the combination of these two experts. As we can see, these are much more accurate experts, nevertheless for all values of α , the combination method improve the objective function - 5.56.% for $\alpha = 0.2$; 6.82% for $\alpha = 0.1$; 13.71% for $\alpha = 0.067$; and 17.20% for $\alpha = 0.05$ compared to the values of the objective function for the better expert.

α	cp			gsc		
	err(%)	rej(%)	opt	err(%)	rej(%)	opt
1/5	0.89	10.96	15.44	1.16	4.99	10.80
1/10	0.78	12.00	19.79	1.00	6.19	16.14
1/15	0.67	13.73	23.84	0.90	7.57	21.14
1/20	0.47	16.25	25.61	0.85	8.58	25.58

Table 3.2.3: Performance of cp and gsc experts

α	err(%)	rej(%)	opt
1/5	1.10	4.73	10.22
1/10	0.82	6.85	15.05
1/15	0.66	8.41	18.25
1/20	0.60	9.23	21.17

Table 3.2.4: Performance of cp-gsc expert combination

3.2.1 Conclusion

We have investigated the simplified version of type 3 multi-expert systems in which each expert outputs a class label and a corresponding confidence. We have developed a general theoretical framework for optimal posterior-probability based combination scheme and have shown that it needs a huge representative training set as well as large memory. This is impractical. We have therefore developed an empirical approach to approximate the joint accuracy distribution function. In this approach, we develop a special measurement, accuracy, which is applicable to all experts. We characterize each expert with a class of accuracy maps. We also develop a family of special combination functions. Finally, we have discussed the optimal accuracy threshold for the recognition decision.

This approach doesn't require mutual independence of experts. In fact, [31] is just a special case of our approach. However, all desirable properties exist from the statistical point of view. A "large enough" and "well represented" training sample set must be available. If only few samples are collected randomly and carelessly, the desired properties of this method cannot be guaranteed [31].

3.3 Fuzzy Control Structure in Multiple Parameter Systems

Most practical software systems are a conglomerate of many modules, where each module has its own parameters that control the accuracy of the module. While each individual module can be optimized by tuning the relevant parameters, it is a non-trivial task to optimize the entire system. When the number of modules and the parameters are few, manual choosing of all parameters is possible by a "trial and error" mechanism. However, when the modules are many, other methods have to be adopted. We use the Fuzzy set technology for this purpose in the application of a Handwritten Address Interpretation System.

We use the Fuzzy set technology for the purpose in the Handwritten Address Interpretation System. The modules that constitute the system are: (i) Address Block Location, (ii) Line separation module, (iii) Word Separation Module, (iv) ZIP Code recognition module, (v) Street number recognition module, and (vi) Street name recognition module. Each of these modules have submodules and face the same problems at a micro-level. The ZIP Code recognition and Street number modules invoke a numeral recognition module and a Street number recognition module invokes a character recognition module. Each of the modules and sub-modules are controlled by many parameters.

In the final analysis, a particular recognition of the address block is deemed either worthy of being accepted as correct or rejected as unreliable. The idea is to accept as high a percentage of addresses as possible while making a minimum of errors. A conservative setting of the thresholds at the final stage would cause many correct results to be rejected in order to keep the error rate low.

We will use a fuzzy measure based on the topology of the "blind" set of handwritten addresses and a pre-determined neighborhood to construct the fuzzy membership function. This will allow a non-linear partitioning of the results to maximize the correct rate. On a test set of 10,000 images, the fuzzy methodology accomplishes 5% higher accuracy.

3.3.1 Handwritten Address Interpretation

Determining the delivery location for mail pieces based on handwritten addresses is a problem that trained humans can normally solve. As a problem in machine reading and interpretation, it presents many challenges like word recognition, digit recognition, combining high level and low level information, control structure to satisfy mutual constraints.

The United States Postal Service sorts mail by utilizing an encoding of each mail piece by a string of digits, the first five of which are referred to as the ZIP Code, and together with the next four as the ZIP+4.

A complete system for unconstrained address interpretation needs several components. At the pre-processing level, the written matter from the handwritten address block image must be distinguished from the background, and irrelevant streaks such as underlines must be removed. The text line locations in the address must be determined, the text lines separated into words, and candidates for the ZIP Code, street number, and street name must be determined. Separating an address line into words is non-trivial because the distance between components is often misleading. The ZIP Code and street number are read next. The digit strings in these fields often have touching digits. Finally, the Postal directory is queried with the street number and

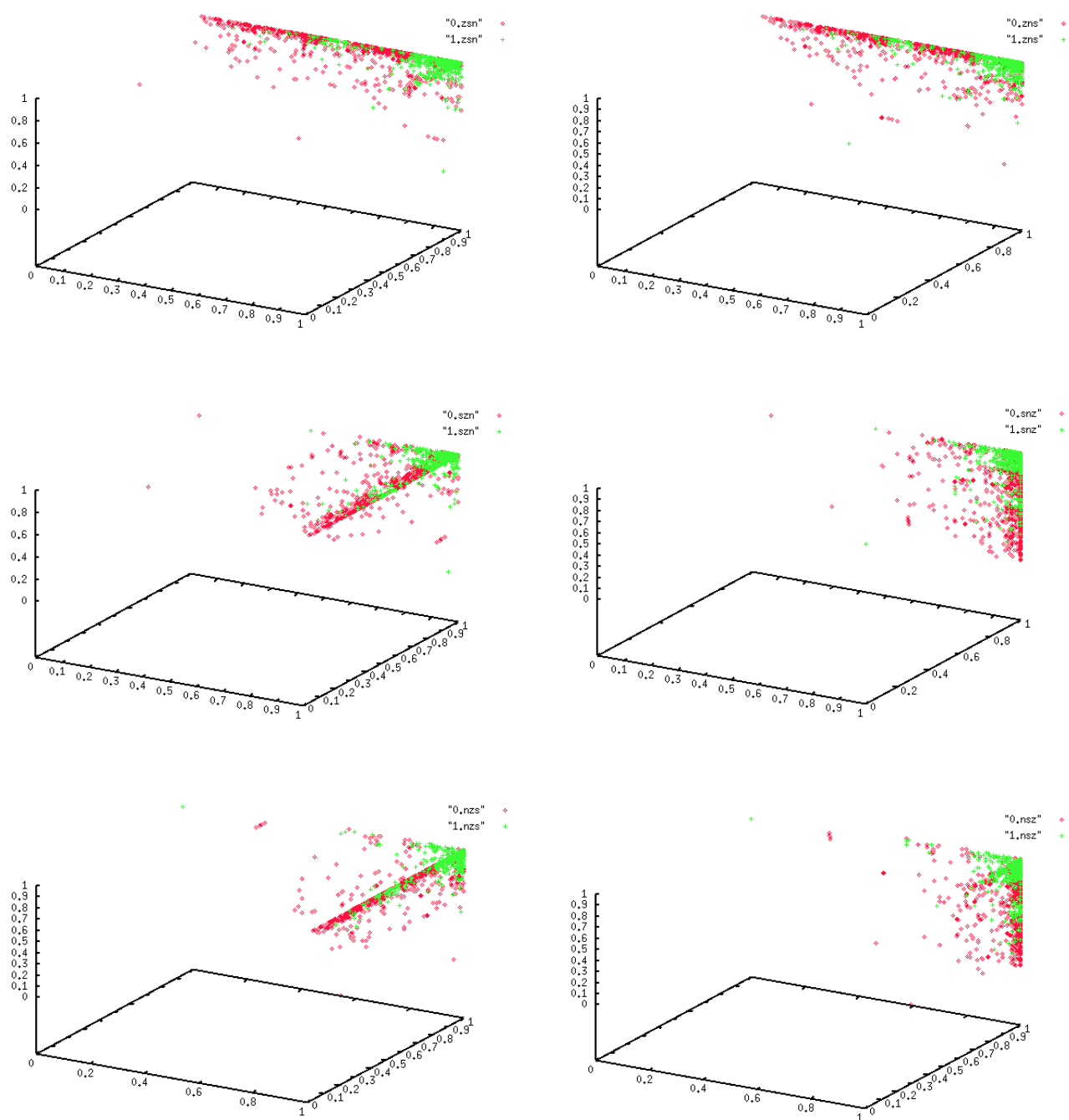
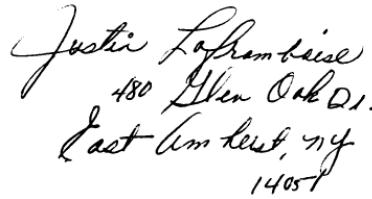


Figure 3.3.1: Figure 1: Each graph shows a perspective of the 3-dimensional space of confidences returned by the ZIP Code recognition module, the Street Number recognition module and the Street Name module. The green points in the space represent addresses that were recognized correctly and the red ones represent the incorrect ones. Clearly a fuzzy boundary exists between the corrects and the errors.

A handwritten address in cursive script. The first line reads "Justin Lofgren". The second line reads "480 Glen Oak Dr.". The third line reads "East Amherst, NY". The fourth line reads "14051".

Encode: 14051-1144

Figure 3.3.2: Figure 2: A handwritten address has many fields that must be recognized in order to find the destination code. The fields of interest in this section are the ZIP Code, the Street Number and the Street Name.

the ZIP Code to obtain a lexicon of street names. The word recognition algorithm ranks the lexicon based on several features it extracts from the street name image. A correct recognition of the ZIP Code AND the Street Number AND the Street Name leads to a correct encoding of the entire address [48].

As is clear from the example of Figure 3.3.2, the accurate location by these modules is a challenge. Recognition of handwritten characters in an unconstrained environment as is presented by the postal domain, adds the difficulty.

The goal is to return an encode for an address only when the system is confident. The cost of rejecting an address image without any encode is less severe than returning an incorrect encode. Hence, each module is expected to return its recognition result along with a confidence measure. During the training phase, the optimal threshold on the confidence measure is chosen so that whenever the confidence exceeds the pre-determined threshold the recognition results of the module is accepted.

3.3.2 Intuitive Understanding of Methodology

The heuristic method uses several thresholds on each confidence measure in order to find the addresses whose recognition results should be accepted. The thresholds themselves are fixed during a training phase which is essentially a “trial and error” procedure. The fuzzy methodology on the other hand partitions the space of all the confidence measures under consideration into a grid. Each grid is labeled as an “Accept” grid or a “Reject” grid in the training phase using a majority rule, *i.e.*, if there are more number of addresses that fall in a grid that are correct, the grid is labeled as an “Accept” grid. The testing phase is a matter of finding the label of the grid a given address falls in.

3.3.3 Topology, Measure and Membership of Fuzzy Sets

Let V be a “blind” set of handwritten addresses, that we will use to simulate the fuzzy boundaries between the two classes ω_0 and ω_1 . V is partitioned into a disjoint union of subsets V_0 and V_1 , where subset V_1 contains all correctly recognized handwritten addresses in V and V_0 is the compliment of V_1 in V , *i.e.* it contains those addresses in V for which at least one of the recognition modules returns result, different from the truth.

Let $f_i(x)$ be the characteristic function of the set V_i , *i.e.*

$$f_i(x) = \begin{cases} 1 & \text{if } x \in V_i \\ 0 & \text{otherwise,} \end{cases}.$$

If x is handwritten address, a confidence vector

$$\Phi(x) = [\phi_1(x), \phi_2(x), \phi_3(x)]^T$$

is formed, where $\phi_1(x)$ is the score (confidence) from the ZIP code recognition, $\phi_2(x)$ is the score (confidence) from the street number recognition, $\phi_3(x)$ is the score (confidence) from the street name recognition.

Let $W = [0, 1]_0 \times [0, 1]_1 \times [0, 1]_2$, be a Cartesian product of 3 copies of the interval $[0, 1]$. Here $[0, 1]_i$ denotes the i^{th} -copy of the interval $[0, 1]$. We can convert W into a topological space by introducing one of the following two topologies: Finite or Euclidean.

1. *Finite Topology:*

Assume that each interval $[0, 1]_i$ is partitioned into finitely many disjoint subintervals $W_i^1, W_i^2, \dots, W_i^m$. Then the set \mathcal{B} of neighborhoods $W_0^{i_0} \times W_1^{i_1} \times W_2^{i_2}$ forms a base of a topology in W . Let us enumerate the elements of \mathcal{B} , i.e $\mathcal{B} = \{B_1, B_2, \dots, B_M\}$, where $M = m^3$.

Let \mathcal{M} be the family of all sets of the form

$$\mathcal{M} = \{B_{i_1}, B_{i_2}, \dots, B_{i_j} : j \in \{1, 2, \dots, M\}, \\ 1 \leq i_1 < i_2 \cdots < i_j \leq M\}.$$

Then \mathcal{M} is a topology of W .

2. *Euclidean Topology:*

In this case we will consider W with the topology induced by the Euclidean topology of the space \mathbb{R}^3 , i.e the topology that is generated by the collection of all balls in the set W [54].

Next step is to construct a σ -algebra, that includes the topology of the set W . In the finite case, it is apparent that \mathcal{M} is a σ -algebra [66].

In the Euclidian case, we define \mathcal{M} to be Borel σ -algebra, generated by the Euclidian topology.

In the finite case, we have a special set of neighborhoods

$$\mathcal{B} = \{B_1, B_2, \dots, B_M\}.$$

Let us notice that for every $y \in W$, there exists a unique neighborhood B_y , containing y . In Euclidian case, given a function $R : W \rightarrow \mathbb{R}$, the special set of neighborhoods includes all balls of the forms

$$B_y = \{z \in W : \text{dist}(y, z) \leq R(y)\}.$$

Now, we will introduce positive σ -additive measures for the σ -algebra \mathcal{M} . Having defined the function $f(x)$, we can define positive σ -additive measures $\mu_i : \mathcal{M} \rightarrow \mathbb{R}$, $i = 0, 1$ as follows:

$$(3.3.1) \quad \mu_i(A) = \sum_{\{x \in V_i : \Phi(x) \in A\}} f_i(x)$$

for all $A \in \mathcal{M}$, as well as the following measure

$$\mu(A) = \mu_0(A) + \mu_1(A).$$

For every $y \in W$, let $\chi_y : W \rightarrow \mathbb{R}$ be the characteristic function of its neighborhood B_y . Let us also define the function

$$\varphi(y) = \int_W \chi_y(z) d\mu(z) = |\{x \in V : \Phi(x) \in B_y\}|,$$

as well as the function

$$\varphi_i(y) = \int_W \chi_y(z) d\mu_i(z) = |\{x \in V_i : \Phi(x) \in B_y\}|,$$

ZIP Code Coordinates		Street Number Coordinates		Street Name Coordinates		Majority		
x_1	x_2	y_1	y_2	z_1	z_2	Correct	Error	Accept / Reject
0.000	0.495	0.830	1.000	0.830	1.000	41	517	Reject
0.495	0.610	0.830	1.000	0.830	1.000	154	17	Accept
0.610	0.830	0.610	0.830	0.830	1.000	4	22	Accept
0.610	0.830	0.830	1.000	0.830	1.000	45	572	Accept
0.830	1.000	0.000	0.495	0.830	1.000	2	2	Reject
0.830	1.000	0.495	0.610	0.830	1.000	2	3	Accept
0.830	1.000	0.610	0.830	0.830	1.000	28	87	Accept
0.830	1.000	0.830	1.000	0.495	0.610	3	4	Accept
0.830	1.000	0.830	1.000	0.610	0.830	17	38	Accept
0.830	1.000	0.830	1.000	0.830	1.000	2016	426	Reject

Table 3.3.1: A grid is defined by the [x,y,z] coordinates. Each grid has addresses inside that are either correct or in error. The majority decides the label of the grid.

for $i = 0, 1$.

Now, we are in a position to define the fuzzy sets corresponding to the different classes. We can define the fuzzy membership functions on the set of all unknown pattern as follows:

$$\xi_i(x) = \begin{cases} \frac{\varphi_i(\Phi(x))}{\varphi(\Phi(x))} & \text{if } \Phi(x) \text{ is in the fuzzy area} \\ \phi_i(x) & \text{otherwise} \end{cases}.$$

Each of these functions $\xi_i(x)$ is the membership function of the corresponding class ω_i . We say that the unknown pattern x belongs to the class ω_i , if $\xi_j(x) < \xi_i(x)$ for all $j \neq i$. Having defined our membership functions $\xi_i(x)$, $i = 0, 1$ we are able to partition the space of all images into class regions and boundaries, that separates them in the following way. The set

$$\Omega_i = \{x : \xi_i(x) > \xi_j(x) \text{ for all } j \neq i\},$$

$i = 0, 1$ is the region, all points of which are in class ω_i .

3.3.4 Experiments

We used a training set of 4,000 handwritten address images. The size of the grids is determined by observing the spread of points in the 3D space. Table 3.3.1 shows the lookup table derived by the method described in this section.

Thresholds			Results	
ZIP Code	Street Number	Street Name	Accept Rate	Correct Rate
0.61	0.61	0.61	79.88	65.27
0.83	0.61	0.61	63.69	79.97
0.83	0.83	0.61	60.73	82.73
0.83	0.83	0.83	59.47	84.03
Fuzzy			64.62	84.73

Table 3.3.2: Results from manual setting of thresholds on the ZIP Code, Street Number, and Street Name yield results which are inferior to the results from the fuzzy methodology.

Table 3.3.2 shows the results on a test set of 10,000 handwritten address images. Of the many modules that are part of a handwritten address interpretation system, we have taken the three most important modules to show the viability of using a fuzzy control instead of a heuristic approach. The three important modules are the ZIP Code recognition module, the Street Number recognition module and the Street name recognition module. First, we manually set the thresholds on the confidence scores returned by the modules as shown in Table 3.3.2. For instance, when all three acceptance thresholds are 0.61, 7,988 addresses are accepted and the remaining are rejected as at least one of the three modules has confidence below the threshold. Of the ones accepted, 6,527 (65.27%) are correct. The fuzzy approach described in this section has an acceptance rate of 64.62% and a correct rate of 84.73%. For an equivalent acceptance rate the heuristic method has a correct rate of 79.97%, almost 5% lower.

3.4 Architecture for Classifier Combination Using Entropy Measure

Machine recognition of isolated handwritten words, especially cursive script, is a difficult problem. The problem is made tractable only when constrained by a *lexicon*, a list of words that includes the truth of the image as one of its elements. Research in handwritten word recognition (HWWR) has traditionally focused on relatively small lexicons, typically comprised of 10 - 1000 entries. Features extracted from the image are matched against every lexicon entry by an expensive matching algorithm, and a confidence value is computed for each lexicon entry [37]. The lexicon entries are ranked in decreasing order of confidence. While this paradigm has proven to be sufficient for many applications from check amounts to street names in mail addresses, there are other applications wherein the lexicons are large (of the order of 10,000 entries or more) and it is no longer practical to compare the features extracted from the image with every lexicon entry. Some means of rapidly eliminating large parts of the lexicon as being unlikely matches is called for. This process is called lexicon reduction, and serves to rapidly trim the original lexicon down to a tractable size for a word classifier.

It is a known fact that classifier performance declines with increasing size of lexicon. This may be attributed to the presence in large lexicons of several entries that the classifier finds difficult to distinguish from the reference. By eliminating some of these entries, lexicon reduction results in improved recognition performance.

Classifiers can combine in either series or in parallel. Figure 3.4.1 shows the architectures possible. This section is about developing a general theory for combination of

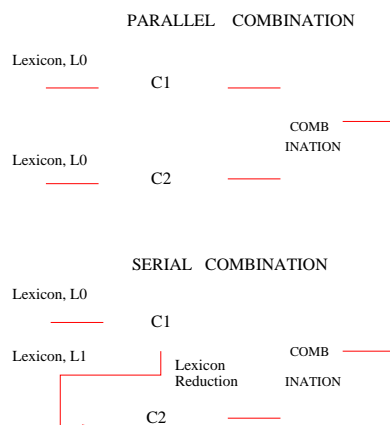


Figure 3.4.1: Two Classifier Combination Models. Parallel Combination methods can have the classifiers acting independently in terms of both looking at the same original lexicon. Serial methods are clearly dependent on each other as the second classifier in the series depends on the results of the first classifier.

Architecture	A_x	A_y	A_1
$2768 \xrightarrow{\text{classifier 1}} 50 \xrightarrow{\text{classifier 2}} 10 \xrightarrow{\text{classifier 3}} 1$	$A_{50} = 70.87\%$	$A_{10} = 39.23\%$	34.40%
$2768 \xrightarrow{\text{classifier 2}} 500 \xrightarrow{\text{classifier 1}} 50 \xrightarrow{\text{classifier 3}} 1$	$A_{500} = 62.77\%$	$A_{50} = 60.77\%$	46.60%
$2768 \xrightarrow{\text{classifier 1}} 50 \xrightarrow{\text{classifier 3}} 1$		$A_{50} = 70.87\%$	52.40%
$2768 \xrightarrow{\text{classifier 1}} 10 \xrightarrow{\text{classifier 3}} 1$		$A_{10} = 66.43\%$	53.83%

Table 3.4.1: There are $n!$ different ways in which n classifiers can be arranged in series. A_x indicates the accuracy of reduction after first stage of classification with classifier1 and A_y after the second stage. x and y are the sizes of the reduced lexicons.

classifiers. In particular, we want to discuss the theoretical and complexity issues pertaining to lexicon reduction and the serial combination of classifiers. Let us assume for the discussion in this section that the word recognizer is the classifier and the lexical entries are the classes. Further, let us begin by discussing the case of two classifiers. We shall later see how the methodology developed can be readily generalized to any number of classifiers. The method of parallel combination of classifiers would submit the same lexicon to both classifiers and combine the results using a variety of methods, such as logistic regression and Borda Count [28],[83]. On the other hand, in serial combination, classifiers that operate later in the engine deal with smaller lexicons. In fact, lexicon reduction is central to the serial combination methods.

Table 3.4 shows a few serial combination methods using 3 classifiers. We take note of the fact that the number of possibilities is very large. First, if there are 3 classifiers, they can be ordered in $3! = 6$ ways. Further, given the original lexicon L_0 of length $|L_0|$, the reduced lexicon output by classifier C_1 can be of $|L_0 - 2|$ different sizes (not counting the cases when the entire lexicon or just 1 entry are returned). For $L_0 = 2768$, it amounts to $6 \times 2766 = 16,596$ possible configurations for the architecture.

The motivation for this section stems from the difficulty confronted by a designer in choosing the correct architecture. Empirical methods are typically used by researchers. However, as the number of classifiers increases and the size of the lexicon (L_0) increases, the situation quickly gets out of hand. Our objective is to develop the theory that will provide the guidelines for choosing the best method of serial combination.

If one were to consider parallel combination methods as well, the possible configurations increase further. Table 3.4 shows how different architectures can be configured

Architecture	A_x	A_y	A_z	A_1
	$A_{500} =$	$A_{50} =$	$A_{10} =$	
	62.77%	26.80%	57.10%	46.83%
	$A_{100} =$	$A_{10} =$	$A_5 =$	
	70.87%	70.87%	27.77%	52.40%

Table 3.4.2: Various configurations are possible when parallel and serial methods are mixed. A_x , A_y , and A_z have the same meaning as above.

mixing the notion of serial and parallel combination. We have shown just a few examples. It should be apparent that the possibilities are very large. We will describe a universal combination architecture that will allow us to enumerate every possible configuration.

Searching for the optimal choice of the architecture is clearly an open research problem. We will develop a greedy algorithm that uses entropy measures to search the optimal architecture. We have experimented with 3 classifiers and present our results (Table 3.4).

3.4.1 Universal Combinator

Let us consider the following model for possible combinations of N classifiers C_1, C_2, \dots, C_N (Figure 3.4.2). Without loss of generality we can assume that the classifiers are running in the same order as they are enumerated. Given an unknown pattern x and a lexicon L_1 , the run of the first classifier C_1 produces a ranked list R_1 of the words in the lexicon and their associated probabilities. A part R_1^0 of that ranked list

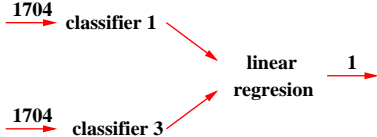
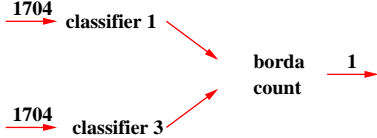
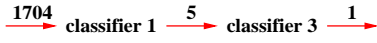
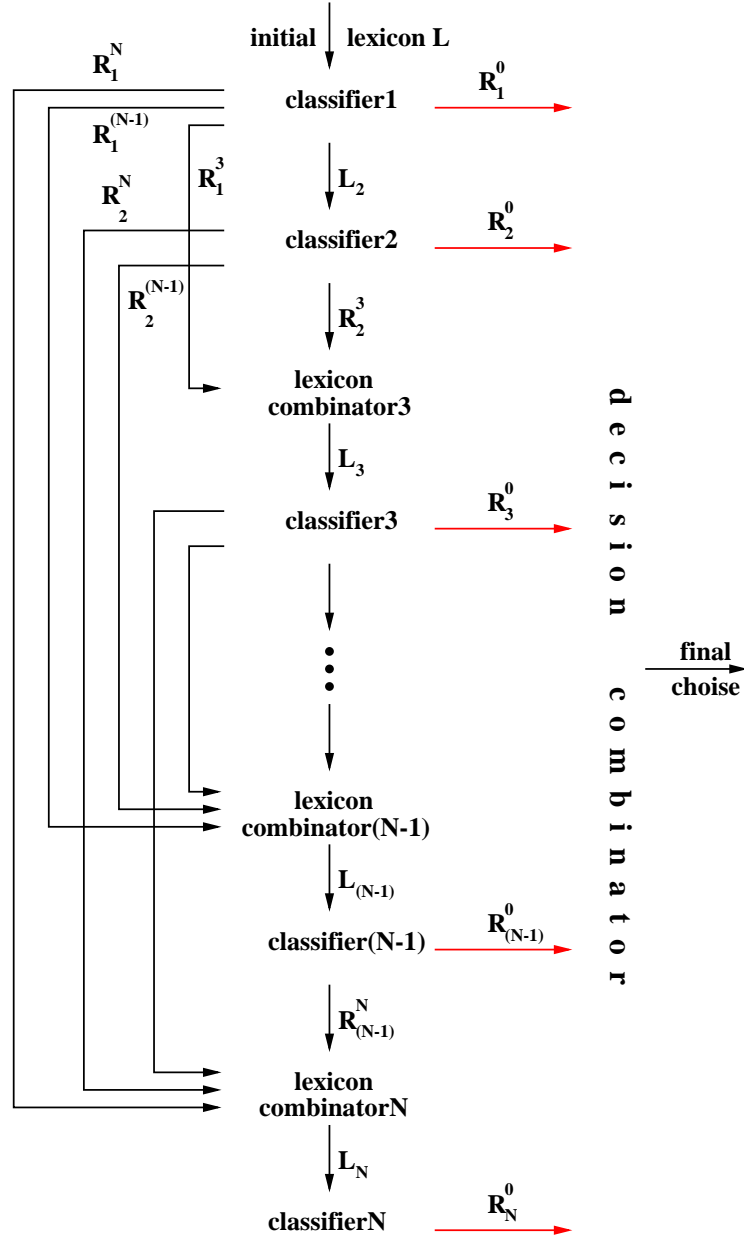
Architecture	A_1
	79.08%
	79.74%
	83.01%
GREEDY (our method described here)	84.31%

Table 3.4.3: The GREEDY algorithm described in this section performs better than the traditionally used parallel combination methods and the *ad hoc* serial combination methods.

is sent to the final decision maker, a part R_1^2 contributes to the lexicon (for a UNION and/or INTERSECTION with other parts of the lexicon) of the second classifier C_2 , while another part R_1^3 is used for building the lexicon of the third classifier C_3 , and so on.

Parts $R_1^2, R_1^3, \dots, R_1^N$ of the ranked list R_1 , that was output from the run of the first classifier, are used in building lexicons for the classifiers, that run after the first one, R_1^2 for C_2 , R_1^3 for C_3 , \dots , R_1^N for C_N . Now, when the run of the first classifier is over, it is the turn of the second classifier C_2 to run with lexicon L_2 built from part R_1^2 of the ranked list, produced by the first classifier C_1 , that is sent to C_2 . The result of the run of the second classifier C_2 , as before, is a ranked list R_2 of the words in the lexicon L_2 and their associated probabilities. A part R_2^0 of that ranked list is sent to the final decision maker, a part R_2^3 contributes to the lexicon of the third classifier C_3 , while another part R_2^4 is used to build the lexicon of the fourth classifier C_4 , and



L_i - denotes the lexicon used by classifier(i)
 R_i^j - denotest the part of rank list, produced by classifier(i), used in forming of lexicon L_j

Figure 3.4.2: Universal Classifier Combination Model

so on. When the run of the second classifier is over, the third classifier C_3 run with lexicon L_3 built from R_1^3 part of the ranked list, produced by the first classifier C_1 , and part R_2^3 of the ranked list, produced by the second classifier C_2 . The output is a ranked list R_3 of the words in the lexicon L_3 and their associated probabilities, parts of which are used to build lexicons for following classifiers. The same procedure is repeated for all classifiers that follow in the recognition engine. The final decision maker outputs the final decision as a ranked list or the top choice.

Our conjecture is that the universal classifier combination model with suitably chosen parameters $R_1^0, R_1^2, R_1^3, \dots, R_1^N, R_2^0, R_2^3, \dots, R_2^N, \dots, R_{N-1}^0, R_{N-1}^N$ represents all possible classifier combination of N classifiers. The model gets its power from the fact that certain R_p^q values can be 0. In fact, if the only non-zero values are accorded to $R_1^0 \dots R_{N-1}^0$, then the classifiers only contribute to the final decision combinator (Figure 3.4.2) and the architecture becomes purely parallel.

Now, our goal will be to see that so called universal classifier combination model with suitably chosen parameters $R_1^0, R_1^2, R_1^3, \dots, R_1^N, R_2^0, R_2^3, \dots, R_2^N, \dots, R_{N-1}^0, R_{N-1}^N$ covers any possible classifier combination of N classifiers. We will partition the set of all classifiers into groups and ordered them in such way that any two classifiers in each group don't communicate each other and only classifiers in groups with smaller number contribute with parts of their lexicon(or produced by their runs with the corresponding lexicons ordered lists) in the assembling of the lexicons of the classifiers in given group. In view of the comments we made, the first group G_1 will consist of only these classifiers, to which the original lexicon was provided. They obviously don't communicate one another,i.e no any of these classifiers is responsible for the lexicons of the rest of the classifiers. The content of the second group G_2 includes only those classifiers which lexicons depend on the parts of the lexicons (or corresponding ordered lists) of classifiers from group G_1 . Then again it is clear that the classifiers

inside this group don't communicate one another (otherwise their lexicon will depend not only on the classifiers from group G_1 but some of the classifiers in group G_2). The same way we form the rest of the groups. Suppose that there are total of m groups. So we can order the classifiers in each group according to their initial order, i.e. if the classifiers C_i and C_j are in the same group, the classifier C_i is ordered ahead of the classifier C_j if and only if $i < j$. In this way a new linear (any two classifiers are comparable) order " $<$ " has been introduced. For the classifiers $C_i \in G_k$ and $C_j \in G_l$, we have $C_i < C_j$ if and only if one of the following two conditions $k < l$ or $k = l$ and $i < j$ is satisfied. Let us notice that with respect to the new order that was introduced, each classifier submit information to classifiers that follow it in this order or to the final decision maker. This shows that if we order the classifiers in our universal model according this new linear order and pass the same information from one classifier to the other (or to the final decision maker) as in the original combination, the results (and the entire performance) of the two combinations will be identical. At the end let us make the following observation. To describe in full our universal classifier model we need to know how the classifiers are ordered. Once the order is established, then we need another $\frac{N(N+1)}{2}$ parameters to fully describe our model, namely the portions of the ranked list R_1 , sent to the rest of the classifiers and to the final decision maker (total of N parameters), the portions of the ranked list R_2 , sent to all classifiers that follow it (i.e. R_2^3, \dots, R_2^N) and to the final decision maker (total of $N-1$ parameters), etc.. Let us make the following observation. We have introduced a linear order in each group of classifiers in order to construct a universal classifier that duplicate the performance of a given combination of N classifiers. If we have introduced different order in each of the groups, we would have ended up with another model that performs exactly the same as the original combination. The reason for such a behavior is the way we set the groups, i.e. classifiers in each group don't communicate one another, they can do classification in parallel and then the

order is not important. Based on this observation, we can conclude there are at least $|G_1| \cdot |G_2| \cdot \dots \cdot |G_m|$ models to represent the same classifier combination, where $|G_i|$ stands for the number of classifiers in group G_i . Example is show in the following figure.

3.4.1.1 Parameters Measuring The Performance of a Black Box Classifier

Again we assume that the result of a run on an unknown pattern x of the given (black box) classifier C with a lexicon L is a ranked list of the words (or part of them) in the lexicon, where each word is associated with a real number as a measure of the membership of the word and the unknown pattern x into the same class of equivalence. Maybe here is the right place to make the following remark. Different classifiers can operate under different equivalence relations. There are some examples:

- two words belong to the same class of equivalence if and only if they have the same length
- two words belong to the same class of equivalence if and only if they have the same characters
- two words belong to the same class of equivalence if and only if a given string is a part of both words

Let us notice, the number of classes of equivalence for different equivalence relations may vary from very small number (words with the same length) to very large one(words having the exactly the same string of characters).

Arguably, the most important characteristic that measures the performance of the classifier is its accuracy. There are two natural ways to measure the accuracy of

the classifier. One way is to use a threshold that separates the good choices in the ranked list from the bad ones. Then the accuracy of the classifier can be described by functions $F : [0, 1] \times [0, 1] \rightarrow [0, 1]$ and $G : [0, 1] \times [0, 1] \rightarrow [0, 1]$ that have the following properties. Given real numbers τ and u between 0 and 1, $F(\tau, u)$ represent on average in percentage how many times the portion of the correct choices, for which the associated probabilities are above the threshold τ , is greater than u . In similar way, given real numbers τ and v between 0 and 1, $G(\tau, v)$ represent on average in percentage how many times the portion of the wrong choices, for which the associated probabilities are above the threshold τ , is greater than v . In particular, in the cases when only there is only one correct choice among the entries in the lexicon, there is more simple description of the accuracy of the classifier using a function $f : [0, 1] \rightarrow [0, 1]$, such that for given threshold τ $f(\tau)$ shows on average in percentage how many times the associated with the correct choice probability is above the threshold τ . Since the size of the lexicon may vary and the probabilities, associated with the words in the lexicon, depend on the size of the lexicon (since they comply with the natural constrain), a better description of the accuracy of the classifier can be given by functions, similar to the ones described above, which involve as a variable the size of the lexicon. Another natural way to measure the accuracy of the classifier is to use fixed number of top choices instead of a threshold. Then the accuracy of the classifier can be described by functions $F : \mathbb{N} \times \mathbb{N} \times [0, 1] \rightarrow [0, 1]$ and $G : \mathbb{N} \times \mathbb{N} \times [0, 1] \rightarrow [0, 1]$ that have the following properties. Given integer numbers N and n and a real number $u \in [0, 1]$, $F(N, n, u)$ represent on average in percentage how many times the portion of the correct choices presented among the top n choices of the classifier running on a lexicon of size N is greater than u . In similar way, given integer numbers N and n and a real number $v \in [0, 1]$, $G(N, n, v)$ represent on average in percentage how many times the portion of the wrong choices presented among the top n choices of the classifier running on a lexicon of size N

	10	100	1000
Top1	95.80	85.86	68.7
Top2	98.26	91.93	80.2
Top20		98.96	
Top50			97.9

Table 3.4.4: HMM performance

	10	100	1000
Top1	92.4	82.5	68.7
Top2	96.1	87.9	74.4
Top5	98.7	92.1	81.1
Top10		94.1	86.0
Top20		96.2	89.2

Table 3.4.5: HGR performance

is greater than v . As we notice before, in the cases when only there is only one correct choice among the entries in the lexicon, there is more simple description of the accuracy of the classifier using a function $f : \mathbb{N} \times \mathbb{N} \rightarrow [0, 1]$, that for given size of the lexicon, used by classifier shows on average in percentage how many times the correct choice is among top n entries of the corresponding ranked list. Another, even more precise description of the accuracy of the classifier can be given by involving the probabilities, associated with top n choices in the following way. We can divide the interval $[0, 1]$ into k (non-intersecting) parts(bins), and we can measure the accuracy of the classifier using the spread of the top n choices over these bins. Let us notice that characteristic may not be very useful in ordering the classifiers, but can be very useful in decision how to split the ranked list of each classifier.

Another also important characteristic of the classifier is the average time $T(m)$ the classifier consume to perform its job on a lexicon of a size m . We can go to a more precise description (in case of word recognizers, using the characters in the lexicon instead of the size).

3.4.1.2 Again about classifier combination

Now, we have the universal model of classifier combination. For any given classifier combination there are an appropriate order of classifiers and suitably chosen parameters, for which the corresponding universal model represents that particular classifier combination. Given N classifiers (as black boxes), the ultimate goal is to find the optimal combination of the classifiers, i.e the appropriate order of the classifiers in the universal model and corresponding parameters. Moreover, different combinations may be suitable for different situations, i.e we may look for a predetermine classifier combination, independent of the pattern or the lexicon, or we may look for a dynamic classifier combination that depend on given pattern and lexicon, or part of them. There are different objective in constructing a classifier combination. Main objective is arguably the accuracy of the combination, but there are other as time consumed for the classification of the process, computational complexity of other resources used.

As we can see from the examples of the about the classifiers characteristic (accuracy), given relatively small lexicon, the classifiers perform well. But if the lexicon size increases, the accuracy decreases significantly and in case of parallel like combinations, the time used for classification in contrast increases significantly. Based on these observations, in this particular situation (paper) our goal is much smaller, namely, given the accuracy rates of each of N classifiers, we want to find a combination of all or part of these classifiers that improves the accuracy rate of any individual classifier, as well to use reasonable time for classification.

Again we have N classifiers C_1, C_2, \dots, C_N . Given a lexicon L_1 and a unknown pattern x , let us look at the entropy in the serial combination at each of the levels (after classification performed by C_1 , after classification performed by C_2 and so on). Let p_w is the probability in the ranked list, produced by the first classifier C_1 ,

associated with the word w from the lexicon. Let L_i be the lexicon that is used by the classifier C_i . Then the initial entropy is given by

$$E_1 = - \sum_{w \in L_1} p_w \ln p_w \quad .$$

Once the lexicon L_1 is split into two, L_2 and $L_1 \setminus L_2$, the new using the Bayesian rule we compute the new conditional probabilities for words in each of both parts. Let H_2 is the entropy associated with the lexicon L_2 and G_2 is the entropy associated with the compliment in L_1 of the lexicon L_2 , i.e $L_1 \setminus L_2$. Then the total entropy is given by

$$E_2 = \alpha_2 H_2 + (1 - \alpha_2) G_2 \quad ,$$

where α is the average in percentage how many times the correct choice is among top $|L_2|$ entries of the ranked list, produced by classifier C_1 using a lexicon of size $|L_1|$. Once the lexicon L_2 is split into two, L_3 and $L_2 \setminus L_3$, the new using the Bayesian rule we compute the new conditional probabilities for words in each of both parts. Let H_3 is the entropy associated with the lexicon L_3 and G_3 is the entropy associated with the compliment in L_2 of the lexicon L_3 , i.e $L_2 \setminus L_3$. Then the total entropy is given by

$$E_3 = \alpha_2 (\alpha_3 H_3 + (1 - \alpha_3) G_3) + (1 - \alpha_2) G_2 \quad ,$$

where α is the average in percentage how many times the correct choice is among top $|L_3|$ entries of the ranked list, produced by classifier C_2 using a lexicon of size $|L_2|$. And we continue in the same way until all classifiers are used.

Now, we will describe a “greedy” algorithm to construct (dynamically) a N classifier combination. Having a pattern x and a lexicon L , we apply first the classifier with best characteristics on lexicons of size $|L|$. Then we will look to cut the lexicon in order to minimize the entropy. Then this reduced lexicon is send to the next classifier. (We have different choices what to send to the other classifiers-we can fix a number(let say

10 or 0 for serial combination), that and top 10 choices are going to be send to the others, but to the second classifier, or maybe we can send the reduced lexicon to all of them). For lexicon combination we use union or intersection of incoming parts. We continue the procedure until we cannot reduce the entropy or use the last classifier. The top choice of the last used classifier is our choice.

3.4.2 Preliminary Results

Here we will give the proofs of few preliminary results that we will need.

Lemma: Suppose that $x, y \geq 0$, such that $x + y = a$, where a is a real number, such that $0 \leq a \leq 1$. Then

$$a \ln \frac{a}{2} \leq x \ln x + y \ln y \leq a \ln a$$

Proof: Let $f(x) = x \ln x + (a - x) \ln(a - x)$.

Then $f'(x) = \ln x + 1 - \ln(a - x) - 1 = \ln \frac{x}{a - x}$.

If $x \in [0, \frac{a}{2})$, then $a - x > \frac{a}{2} > x$, so $\frac{x}{a - x} < 1$ and $\ln \frac{x}{a - x} < 0$,

i.e $f'(x) < 0$ for $x \in [0, \frac{a}{2})$. Therefore, $f(x)$ is decreasing in $[0, \frac{a}{2})$.

If $x \in (\frac{a}{2}, a]$, then $a - x < \frac{a}{2} < x$, so $\frac{x}{a - x} > 1$ and $\ln \frac{x}{a - x} > 0$,

i.e $f'(x) > 0$ for $x \in (\frac{a}{2}, a]$. Therefore, $f(x)$ is increasing in $(\frac{a}{2}, a]$.

Then, $f(\frac{a}{2}) \leq f(x) \leq \max(f(0), f(a))$ for $x \in [0, a]$.

$$f(\frac{a}{2}) = \frac{a}{2} \ln \frac{a}{2} + \frac{a}{2} \ln \frac{a}{2} = a \ln \frac{a}{2}$$

$f(0) = 0 \ln 0 + a \ln a = a \ln a$, $f(a) = a \ln a + 0 \ln 0 = a \ln a$. Then,

$$a \ln \frac{a}{2} \leq x \ln x + y \ln y \leq a \ln a$$

Corollary 1:

$$(x + y) \ln \frac{x + y}{2} \leq x \ln x + y \ln y \leq (x + y) \ln(x + y)$$

Corollary 2:

$$\text{A. } x_1 \ln x_1 + x_2 \ln x_2 + \cdots + x_n \ln x_n \leq (x_1 + x_2 + \cdots + x_n) \ln(x_1 + x_2 + \cdots + x_n)$$

$$\text{B. } x_1 \ln x_1 + x_2 \ln x_2 + \cdots + x_n \ln x_n \geq (x_1 + x_2 + \cdots + x_n) \ln \frac{x_1 + x_2 + \cdots + x_n}{n}$$

Proof: A. By the corollary 1,

$$\begin{aligned} x_1 \ln x_1 + x_2 \ln x_2 + \cdots + x_n \ln x_n &\leq \\ &\leq (x_1 + x_2) \ln(x_1 + x_2) + \cdots + x_n \ln x_n \leq \\ &\leq (x_1 + x_2 + x_3) \ln(x_1 + x_2 + x_3) + \cdots + x_n \ln x_n \leq \\ &\leq \cdots \leq (x_1 + x_2 + \cdots + x_n) \ln(x_1 + x_2 + \cdots + x_n) \end{aligned}$$

B. The set $S_c = \{(x_1, x_2, \dots, x_n) : x_i \geq 0 \text{ for } i = 1, 2, \dots, n \text{ and } x_1 + x_2 + \cdots + x_n = c\}$ is compact and function $f(x_1, x_2, \dots, x_n) = x_1 \ln x_1 + x_2 \ln x_2 + \cdots + x_n \ln x_n$ is continuous and then it reaches its minimum value in S_c . Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is the point, where $f(x_1, x_2, \dots, x_n)$ reaches its minimum value. Then, let us suppose that $X_i \neq X_j$. Let

$$Y_k = \begin{cases} X_k & \text{if } k \neq i, j \\ \frac{X_i + X_j}{2} & \text{if } k = i, j \end{cases}.$$

Obviously, $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n) \in S_c$. Then,

$$\begin{aligned} f(X_1, X_2, \dots, X_n) - f(Y_1, Y_2, \dots, Y_n) &= \\ &= X_i \ln X_i + X_j \ln X_j - (X_i + X_j) \ln \frac{X_i + X_j}{2} > 0, \end{aligned}$$

i.e $f(X_1, X_2, \dots, X_n) > f(Y_1, Y_2, \dots, Y_n)$, which contradicts with the minimum of f at \mathbf{X} . So, $X_i = X_j$ and then

$$\begin{aligned} f(x_1, x_2, \dots, x_n) &\geq \frac{x_1 + x_2 + \dots + x_n}{n} \ln \frac{x_1 + x_2 + \dots + x_n}{n} + \\ &\frac{x_1 + x_2 + \dots + x_n}{n} \ln \frac{x_1 + x_2 + \dots + x_n}{n} + \dots + \frac{x_1 + x_2 + \dots + x_n}{n} \ln \frac{x_1 + x_2 + \dots + x_n}{n} = \\ &= (x_1 + x_2 + \dots + x_n) \ln \frac{x_1 + x_2 + \dots + x_n}{n}. \end{aligned}$$

3.4.3 Entropy Measure for Serial Classifier Combination

Let C, C_1, C_2, \dots, C_m are classifiers, that participate in a serial classifier combination. Without loss of generality, we can assume that C is the first one in the combination. Given a lexicon L , an unknown pattern x , and a classifier C which assigns a probability p_w to each word (w) in the lexicon, the initial entropy of the system is given by

$$E = - \sum_{w \in L} p_w \ln p_w \quad .$$

Our conjecture is that (under certain conditions) the entropy monotonically reduces as the lexicon keeps getting smaller. There are two cases to be considered. We will show

here that the conjecture holds if the classifiers are error free. That is at every stage, the classifier preserves the true choice in the reduced lexicon. Since the last classifier in a serial engine returns just 1 choice, this case assumes that the final recognition choice returned by the cascade of classifiers is correct. In fact, the conjecture holds for cases when the classifiers are not necessarily error-free.

Lexicon L is split into L_1 and $L \setminus L_1$. H_1 is the entropy associated with the lexicon L_1 (with the corresponding *a posteriori* probabilities) and G_1 is the entropy associated with the complement $L \setminus L_1$. The total entropy of the system is given by

$$E_1 = \alpha_1 H_1 + (1 - \alpha_1) G_1 ,$$

where α_1 is a parameter.

Lexicon L_1 is split into L_2 and $L_1 \setminus L_2$. H_2 is the entropy associated with the lexicon L_2 and G_2 is the entropy associated with the complement $L_1 \setminus L_2$. The total entropy of the system is given by

$$E_2 = \alpha_1(\alpha_2 H_2 + (1 - \alpha_2) G_2) + (1 - \alpha_1) G_1 ,$$

where α_1 and α_2 are parameters.

The entropy of the system keeps decreasing as the lexicon gets reduced (under certain conditions). We illustrate the entropy decrease in a simplified situation (the general case is more complicated but it can be treated similarly). Let x be a unknown pattern and L is the corresponding lexicon, that contains the words c_1, c_2, \dots, c_n . The classifiers C_1, C_2, \dots, C_m detects features F_1, F_2, \dots, F_m respectively. Let ε_i be error on average that classifier C_i makes detecting the corresponding feature F_i in patterns that don't have it and let ϵ_i be error on average that classifier C_i makes not detecting the corresponding feature F_i in patterns that do have it. Let S_i be the set of all words from the lexicon L , that have the feature F_i and let $|S_i|$ be its cardinality,

i.e the number of elements of S_i . Let A_i be the event “ F_i is detected” and $A_{i,j}$ be the event “ F_i is detected first and then F_j is detected”. Then obviously $A_{i,j} = A_j \circ A_i$

Remark: If the recognition is error free, i.e $\varepsilon_i = \epsilon_i = 0$ for $i = 1, 2, \dots, m$, then the result of the event $A_{i,j}$ is identical with the result of the event $A_{j,i}$. Indeed, the result of both event is a pattern that has both features F_i and F_j , so the order of classifiers C_1, C_2, \dots, C_m in the serial engine is not important.

Let $Prob(x = c_k) = p_k$ according to classifier C . Then

$$Prob(F_i \text{ is detected in } x | x = c_k) = \begin{cases} \epsilon_i & c_k \notin S_i \\ 1 - \varepsilon_i & \text{otherwise} \end{cases}.$$

$$Prob(F_i \text{ is detected in } x) =$$

$$\begin{aligned} &= \sum_{c_k} Prob(F_i \text{ is detected in } x | x = c_k) Prob(x = c_k) = \\ &= \left(\sum_{c_k \in S_i} + \sum_{c_k \notin S_i} \right) Prob(F_i \text{ is detected in } x | x = c_k) Prob(x = c_k) = \end{aligned}$$

$$= \sum_{c_k \in S_i} (1 - \varepsilon_i) p_k + \sum_{c_k \notin S_i} \epsilon_i p_k = (1 - \varepsilon_i) \sum_{c_k \in S_i} p_k + \epsilon_i \sum_{c_k \notin S_i} p_k$$

$$Prob(x = c_k | F_i \text{ is detected in } x) =$$

$$\frac{Prob(F_i \text{ is detected in } x | x = c_k) Prob(x = c_k)}{Prob(F_i \text{ is detected in } x)} =$$

$$= \begin{cases} \frac{(1 - \varepsilon_i) p_k}{(1 - \varepsilon_i) \sum_{c_j \in S_i} p_j + \varepsilon_i \sum_{c_j \notin S_i} p_j} & c_k \in S_i \\ \frac{\varepsilon_i p_k}{(1 - \varepsilon_i) \sum_{c_j \in S_i} p_j + \varepsilon_i \sum_{c_j \notin S_i} p_j} & c_k \notin S_i \end{cases}.$$

Let $P_i = \sum_{c_j \in S_i} p_j$, $Q_i = \sum_{c_j \notin S_i} p_j$, and $R_i = (1 - \varepsilon_i)P + \varepsilon_i Q$ for $i = 1, 2, \dots, m$. Obviously, $P_i + Q_i = 1$. for for $i = 1, 2, \dots, m$. For simplicity, we set $P = P_1$, $Q = Q_1$ and $R = R_1$, as well as $\varepsilon = \varepsilon_1$, $\epsilon = \epsilon_1$ and $k = |S_1|$.

If after application of classifier C_1 the lexicon L is reduced to the set S_1 (the set of all words that have feature F_1) in the case feature F_1 has been detected or $L \setminus S_1$ in the case feature F_1 has not been detected. The a *posteriori* probability of x to be the word c_k in S_1 is

$$\bar{p}_k = \frac{(1 - \varepsilon_1) p_k}{(1 - \varepsilon_1) \sum_{c_j \in S_1} p_j + \varepsilon_1 \sum_{c_j \notin S_1} p_j}.$$

The a *posteriori* probability of x to be the word c_k in $L \setminus S_1$ is

$$\bar{p}_k = \frac{\varepsilon_1 p_k}{(1 - \varepsilon_1) \sum_{c_j \in S_1} p_j + \varepsilon_1 \sum_{c_j \notin S_1} p_j}.$$

The initial entropy of the system is

$$E_0 = - \sum_{i=1}^n p_i \ln(p_i).$$

The entropy of the set S_1 is

$$\begin{aligned}
H_1 &= - \sum_{c_k \in S_1} \frac{(1 - \varepsilon_1)p_k}{(1 - \varepsilon_1) \sum_{c_k \in S_1} p_k + \epsilon_1 \sum_{c_k \notin S_1} p_k} \ln \frac{(1 - \varepsilon_1)p_k}{(1 - \varepsilon_1) \sum_{c_k \in S_1} p_k + \epsilon_1 \sum_{c_k \notin S_1} p_k} = \\
&= - \frac{1 - \varepsilon}{R} \sum_{c_k \in S_1} p_k (\ln(1 - \varepsilon) + \ln p_k - \ln R) = \\
&= - \frac{1 - \varepsilon}{R} \left(\sum_{c_k \in S_1} p_k \ln p_k + P \ln(1 - \varepsilon) - P \ln R \right).
\end{aligned}$$

The entropy of the set $L \setminus S_1$ is

$$\begin{aligned}
G_1 &= - \sum_{c_k \notin S_1} \frac{\epsilon_1 p_k}{(1 - \varepsilon_1) \sum_{c_k \in S_1} p_k + \epsilon_1 \sum_{c_k \notin S_1} p_k} \ln \frac{\epsilon_1 p_k}{(1 - \varepsilon_1) \sum_{c_k \in S_1} p_k + \epsilon_1 \sum_{c_k \notin S_1} p_k} = \\
&= - \frac{\epsilon}{R} \sum_{c_k \notin S_1} p_k (\ln \epsilon + \ln p_k - \ln R) = \\
&= - \frac{\epsilon}{R} \left(\sum_{c_k \notin S_1} p_k \ln p_k + Q \ln \epsilon - Q \ln R \right).
\end{aligned}$$

The entropy of the system after the reduction is

$$\begin{aligned}
E_1 &= \alpha H_1 + (1 - \alpha) G_1 = \\
&= -\alpha \frac{1 - \varepsilon}{R} \left(\sum_{c_k \in S_1} p_k \ln p_k + P \ln(1 - \varepsilon) - P \ln R \right) -
\end{aligned}$$

$$-(1-\alpha)\frac{\epsilon}{R}\left(\sum_{c_k \notin S_1} p_k \ln p_k + Q \ln \epsilon - Q \ln R\right).$$

Using the preliminary results we get

$$\begin{aligned} E_0 - E_1 &= -\sum_{i=1}^n p_i \ln(p_i) + \\ &+ \alpha \frac{1-\epsilon}{R} \left(\sum_{c_k \in S_1} p_k \ln p_k + P \ln(1-\epsilon) - P \ln R \right) + \\ &+ (1-\alpha) \frac{\epsilon}{R} \left(\sum_{c_k \notin S_1} p_k \ln p_k + Q \ln \epsilon - Q \ln R \right) = \\ &= \left(\alpha \frac{1-\epsilon}{R} - 1 \right) \sum_{c_k \in S_1} p_k \ln p_k - \left(1 - (1-\alpha) \frac{\epsilon}{R} \right) \sum_{c_k \notin S_1} p_k \ln p_k + \\ &+ \alpha \frac{P}{R} (1-\epsilon) \ln(1-\epsilon) + (1-\alpha) \frac{Q}{R} \epsilon \ln \epsilon - \\ &- \alpha (1-\epsilon) \frac{P}{R} \ln R - (1-\alpha) \frac{Q}{R} \ln R \geq \\ &\geq \left(\alpha \frac{1-\epsilon}{R} - 1 \right) P \ln \left(\frac{P}{|S_1|} \right) - \left(1 - (1-\alpha) \frac{\epsilon}{R} \right) Q \ln Q + \\ &+ \alpha \frac{P}{R} (1-\epsilon) \ln(1-\epsilon) + (1-\alpha) \frac{Q}{R} \epsilon \ln \epsilon - \end{aligned}$$

$$\begin{aligned}
& -\alpha(1-\varepsilon)\frac{P}{R}\ln R - (1-\alpha)\frac{Q}{R}\ln R = \\
& = \frac{1}{R}[(\alpha(1-\varepsilon) - R)P\ln\left(\frac{P}{k}\right) - (R - (1-\alpha)\epsilon)Q\ln Q + \\
& + \alpha P(1-\varepsilon)\ln(1-\varepsilon) + (1-\alpha)Q\epsilon\ln\epsilon - \alpha(1-\varepsilon)P\ln R - (1-\alpha)\epsilon Q\ln R].
\end{aligned}$$

Let us define for $y \in [0, 1]$ the following function

$$\begin{aligned}
f(y) = & (\alpha(1-\varepsilon) - (1-\varepsilon)x - \epsilon(1-x))(x\ln x - x\ln m) - \\
& - ((1-\varepsilon)x + \epsilon(1-x))(1-x)\ln(1-x) + \alpha(1-\varepsilon)x + (1-\alpha)\epsilon\ln\epsilon(1-x) - \\
& - \alpha(1-\varepsilon)x\ln((1-\varepsilon)x + \epsilon(1-x)) - (1-\alpha)\epsilon(1-x)\ln((1-\varepsilon)x + \epsilon(1-x)).
\end{aligned}$$

In the case when classifier is error-free and feature F_1 is detected (i.e $\varepsilon = \epsilon = 0$ and $\alpha = 1$) the function $f(x)$ looks like this (Figure 3.4.3).

$$\begin{aligned}
f(y) & = -y\ln y - (1-y)\ln k - (1-y)\ln(1-y) \\
f'(y) & = \ln \frac{k(1-y)}{y}
\end{aligned}$$

Note that if $S \geq r_k$, then $E_0 \geq E_1$, i.e entropy is decreasing

3.4.3.1 Probability Values for Lexical Entries

Most methods take the image of a handwritten word and a lexicon of possible words, and rank the lexicon based on the “goodness” of match between each lexicon entry and

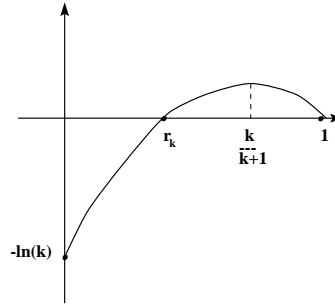


Figure 3.4.3: Graph of $f(y)$ in the case classifier C_1 is error-free

the word image. Typically, the word recognizer computes a measure of “similarity” between each lexicon entry and the word image and uses this measure to sort the lexicon in descending order of the similarity measure. The lexicon entry with the highest similarity is the top choice of the recognizer. The top m choices are often referred to as the confusion set as it contains the lexicon entries that are “similar” to the actual lexicon entry that matches the truth in some feature space.

We have developed elsewhere [7] the groundwork for the use of Bayesian methodology in integration of recognizers with any subsequent processing by deriving meaningful probabilistic measures from recognizers. This allows us to compute the entropy values.

3.4.4 Greedy Algorithm

We have developed a “greedy” algorithm to dynamically construct a combinator of N classifiers. Given a test pattern x and a lexicon L , we first apply the classifier with the best recognition rate on lexicons of size $|L|$. We reduce the lexicon in order to minimize the entropy. This reduced lexicon is sent to the next classifier. We continue the process until we cannot reduce the entropy any more or the last classifier is exhausted. The top choice of the last classifier used is the final recognition choice.

Algorithm

- input: pattern x , lexicon L ,
- for ($i = 1$; $i \leq N$; $i++$) {
 1. choose from unused classifiers the one with the best accuracy performance on a lexicon of size $|L|$
 2. perform the classification and create a ranked list R from the words in the lexicon L according to their confidences
 3. if ($i == N$) {
 - return the first entry of R
 - exit }
 4. for ($k = 1$; $k \leq |L|$; $k++$) {
 - R_1 is formed from the first k words of the ranked list R with new *a posteriori* probabilities
 - R_2 is formed from the rest of the words of the ranked list R with new *a posteriori* probabilities
 - choose appropriate α
(for example: probability of the true choice to be present in the lexicon R_1)
 - compute the entropy H_k }
 5. find the minimal entropy among H_k , $2 < k < |L|$
 6. if ($H_{|L|}$ is the minimum) {
 - mark the classifier as used, but do not use it in the combination
 - keep the same lexicon L }

7. if (H_1 is the minimum) {
 - return the first entry of R
 - exit }
8. if (H_m is the minimum) {
 - mark the classifier as used, but do not use it in the combination
 - keep in the lexicon L only the first m words of the ranked list R

3.4.5 Experiments

Three word classifiers - CMWR (Character Model Word Recognizer) [16] is Classifier₁, WMWR (Word Model Word Recognizer)[37] is classifier₂ and HOL (Holistic) [49] is classifier₃ in our experiments. Each of them takes as input a binary image and an ASCII lexicon, computes a probability each lexicon entry and ranks the lexicon by decreasing probabilities.

WMR is a fast lexicon-driven analytical classifier that operates on the chain-coded description of the street name image. Following slant normalization and smoothing, the image is “over-segmented” at likely character segmentation points. The resulting segments are grouped and the extracted features matched against letters in each lexicon entry using a dynamic programming algorithm.

CMR adopts a different approach. After preprocessing and over-segmentation, segments are grouped in various ways and OCR is performed on the groups to obtain a graph of possible character candidates. For each lexicon entry, the best path through the graph is then determined. CMR is computationally more expensive than WMR.

The two recognizers are sufficiently orthogonal in approach as well as features used to be useful in a combination strategy.

HOL does not perform any segmentation, but uses holistic information such as length, ascenders and descenders to classify the image.

The individual performances of CMR and WMR are shown in Figure 3.4.4. The Oracle represents the method of combination where a correct result is obtained in the top choice if either of the recognizers has it correct. Tables 3.4.6&3.4.7 show how

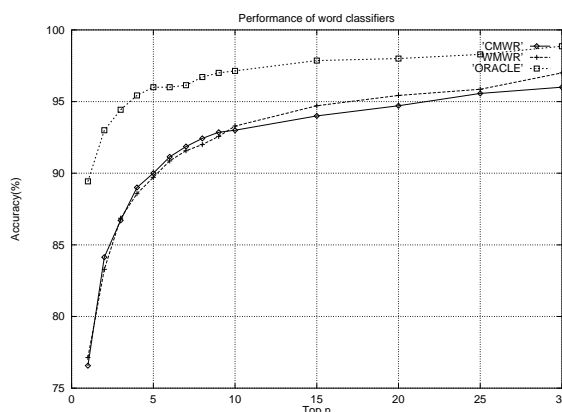


Figure 3.4.4: Oracle combination shows the best possible results that one could obtain from the combination of WMR and CMR

the entropy of the system reduces with lexicon reduction. In this case the classifiers are not error free. However, the entropy still goes down. Table 3.4.8 shows that the GREEDY method finds the lowest entropy when compared to other *ad hoc* methods of determining the architecture of combination.

Image Samples	$\xrightarrow{1704}$ classifier 1 $\xrightarrow{10}$ classifier 2 $\xrightarrow{5}$ classifier 3 $\xrightarrow{1}$			$\xrightarrow{1704}$ classifier 1 $\xrightarrow{50}$ classifier 2 $\xrightarrow{10}$ classifier 3 $\xrightarrow{1}$		
	after C_1	after C_2	final entropy	after C_1	after C_2	final entropy
1	5.577219	3.614541	2.814851	6.294333	4.193750	3.366025
2	5.265298	3.447350	2.646646	5.946801	4.045334	3.217013
3	5.537880	3.594105	2.793283	6.250293	4.175640	3.347035
4	5.047477	3.328180	2.530059	5.701441	3.938757	3.112000
5	5.549012	3.598829	2.799862	6.263592	4.180160	3.352859
6	5.577230	3.615262	2.814670	6.294185	4.194074	3.365882

Table 3.4.6: Note the reduction in entropy after each stage of lexicon reduction on 6 example image samples for two different architectures.

Image Samples	$\xrightarrow{1704}$ classifier 1 $\xrightarrow{5}$ classifier 3 $\xrightarrow{1}$		$\xrightarrow{1704}$ classifier 1 $\xrightarrow{10}$ classifier 3 $\xrightarrow{1}$	
	after C_1	final entropy	after C_1	final entropy
1	3.730102	2.761246	4.064764	2.783761
2	3.554187	2.584103	3.897464	2.615539
3	3.708750	2.738524	4.044601	2.762236
4	3.428433	2.461479	3.778463	2.498956
5	3.713430	2.745451	4.049117	2.768770
6	3.730988	2.761039	4.065277	2.783551

Table 3.4.7: Note the reduction in entropy after each stage of lexicon reduction on 6 example image samples for two different architectures.

Image Samples	"Greedy"	$\xrightarrow{1704} \text{classifier 1} \xrightarrow{5} \text{classifier 3} \xrightarrow{1}$	$\xrightarrow{1704} \text{classifier 1} \xrightarrow{10} \text{classifier 3} \xrightarrow{1}$	$\xrightarrow{1704} \text{classifier 1} \xrightarrow{10} \text{classifier 2} \xrightarrow{5} \text{classifier 3} \xrightarrow{1}$
1	2.76	2.76	2.78	2.81
2	2.58	2.58	2.62	2.65
3	2.74	2.74	2.76	2.79
4	2.46	2.46	2.50	2.53
5	2.75	2.75	2.77	2.80
6	2.76	2.76	2.78	2.81
7	2.65	2.66	2.69	2.72

Table 3.4.8: Final entropy on various image samples using the GREEDY method and other *ad hoc* architectures. Note that the GREEDY method usually has the final lowest entropy.

3.4.6 Remarks

Research in handwritten word recognition has traditionally concentrated on small lexicons of 10 - 1000 words. Several real-world applications, such as the recognition of English prose, involve large lexicons of 10,000 - 50,000 words. Existing classifiers may still be used for these tasks if preceded by a *lexicon reduction* step. The task of lexicon reduction is that of rapidly discarding from the original lexicon entries that are unlikely to match the given image. The resulting two-stage architecture is a *serial combination* or *cascading* of classifiers, and is an effective method of dealing with large lexicons in real-life word recognition scenarios. Moreover, by discarding entries that may potentially confuse the classifier, lexicon reduction results in improved recognition performance. In this section, we have presented an overview of lexicon reduction as a problem in its own right, and discussed some of the issues relating to the design and construction of different combination methods. We have presented a method of serial combination of classifiers using a GREEDY algorithm that strives for minimal entropy of the system. Its recognition accuracy is superior to other *ad hoc* combination methods.

We have shown theoretically that if the classifiers are error free, the entropy of the system must reduce as the lexicon size keeps reducing.

We have introduced the notion of the universal combinator.

Chapter 4

Conclusions

We have contributed to the field of classifier combination on two fronts

- (i) Measures of recognition and
- (ii) Applications of combination

1. Measures:

- (a) We have presented a methodology for deriving pseudo-probability values given recognition scores assigned to digits by a recognizer of type III. We have argued that these pseudo-probability values improve the performance of the recognizer and they also are more informative and useful in a variety of applications including, but not limited to, multiple classifier methodologies.

The results show that the extent of improvement depends on the distribution of the different classes. Also, if the original recognition system performs more “poorly”, DPPS improves the performance of the system by a greater degree.

Our hypothesis is that the performance of the algorithm will improve significantly provided the training set captures the properties of the recognition system, as well as optimal partition limits, needed for the quantization of the validation set. In some cases, it can be proved theoretically that proper validation set in combination with proper partition can improve the recognition system to the theoretically optimal.

DPPS modified classifier is almost as time consuming as the classifier itself in cases when d-neighborhood approach is not used. In the cases, when we choose to use d-neighborhood DPPS, the better correctness rate is a trade-off for increased classification time.

- (b) We have developed a new fuzzy measure for the performance of any given expert that captures the topology of the space formed by scores that the expert returns as output. The benefits of this new fuzzy measure can be considered in two directions: deriving new fuzzy scores as a uniform measure, that can be used in combination of expert systems and improvement of the performance of expert systems.

We have proposed the following two-pronged approach. For the case where a sample falls in “close” proximity of a cluster center (of a class to which it will be assigned), use the traditional method of distances to assign the class and compute its confidence score. For the case where the distances from several cluster centers are all approximately the same, we will adopt a modification of the KNN rule. Such an approach addresses both the observations made above, and at the same time does not cause prohibitive computational issues. Since we expect only a small portion of the feature space to be near the boundaries between classes, the size of the lookup table will be reasonable.

The fuzzy measure is based on the topology of the "blind" set of training samples and pre-determined neighborhoods to construct the fuzzy membership function. This allows a non-linear partitioning of the space, formed by scores that the expert returns as outputs. As a result of this new nonlinear partition of the score space the performance of a given expert is improved. Also we get a new uniform measure that can be used in combination of different types of experts.

2. Applications:

- (a) We have developed a general theoretical framework for optimal posterior-probability based combination scheme and have shown that it needs a huge representative training set as well as large memory. We have proposed the simplified version of type 3 multi-expert systems in which each expert outputs a class label and a corresponding confidence. We have developed an empirical approach to approximate the joint accuracy distribution function. We develop a special measurement, accuracy, which is applicable to all experts. We characterize each expert with a class of accuracy maps. We also develop a family of special combination functions. Finally, we have discussed the optimal accuracy threshold for the recognition decision.

This approach doesn't require mutual independence of experts. However, all desirable properties exist from the statistical point of view. A "large enough" and "well represented" training sample set must be available.

- (b) We have used a fuzzy measure based on the topology of the "blind" set of handwritten addresses and pre-determined neighborhoods to construct the fuzzy membership functions.

The proposed method is very effective, especially in cases when we have to combine the results of modules that use physically different set of features/measurements.

- (c) We have presented an overview of lexicon reduction as a problem in its own right, and discussed some of the issues relating to the design and construction of different combination methods. We have presented a method of serial combination of classifiers using a GREEDY algorithm that strives for minimal entropy of the system. Its recognition accuracy is superior to other *ad hoc* combination methods.

We have shown theoretically that if the classifiers are error free, the entropy of the system must reduce as the lexicon size keeps reducing.

We have introduced the notion of the universal combinator.

Bibliography

- [1] K. Al-ghoneim and B. V. V. Kumar. "Unified decision combination framework". *Pattern Recognition*, 31(12):2077–2089, 1998.
- [2] P. Arabie, L. J. Hubert, and G. D. Soete, editors. "*Clustering and Classification*". World Scientific, River Edge, New Jersey, 1998.
- [3] S. C. Bagui and N. R. Pal. "A multi-stage generation of the rank nearest neighbour classification rule". *Pattern Recognition Letters*, 16(6):601–614, 1995.
- [4] J. A. Barnett. "Computational methods for a mathematical theory of evidence". In *Seventh International Conference on Artificial Intelligence*, pages 868–875, Vancouver, Canada, 1981.
- [5] D. Black. "*The theory of committees and elections*". Cambridge University Press, 1963.
- [6] D. Bouchaffra, V. Govindaraju, and S. Srihari. "A methodology for determining probability of correctness of recognizer scores". In *IEEE Conference on Computer Vision and Pattern Recognition*, Santa Barbara, California, 1998.
- [7] D. Bouchaffra, V. Govindaraju, and S. Srihari. "A methodology for mapping scores to probabilities. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(9):923–927, 1999.

- [8] G. Carapenter, S. Grossberg, N. Reynolds, and D. Rosen. "Fuzzy artmap: a neural network architecture for incremental supervised learning of analog multi-dimensional maps". *IEEE Transactions on Neural Networks*, 3(5):698–713, 1992.
- [9] S. B. Cho and J. H. Kim. "Combining multiple neural networks by fuzzy integral for robust classification". *IEEE Transactions on Systems, Man, and Cybernetics*, 25(2):380–384, 1995.
- [10] S. B. Cho and J. H. Kim. "Multiple network fusion using fuzzy logic". *IEEE Transaction on Neural Networks*, 6(2):497–501, 1995.
- [11] K. Chou, L. Tu, and I. Shyu. "Performance analysis of a multiple classifiers system for recognition of totally unconstrained handwritten numerals". In *Fifth International Workshop on Frontiers in Handwriting Recognition*, pages 480–487, Taiwan, 1994.
- [12] B. V. Dasarathy. "*Nearest neighbor (NN) norms : nn pattern classification techniques*". IEEE Computer Society Press, Los Alamitos, California, 1991.
- [13] R. O. Duda, P. E. Hart, and D. G. Stork. "*Pattern Classification*", 2nd edition. 1997.
- [14] B. Efron. "The efficiency of logistic regression compared to normal discriminant analysis". *Journal of American Statistical Association*, 70(5):892–898, 1975.
- [15] A. M. Erkmen and H. E. Stephanou. "Informational fractals for evidential pattern classification". *IEEE Transactions on Systems, Man, and Cybernetics*, 20(5):1103–1114, 1990.
- [16] J. T. Favata and S. N. Srihari. "Off-line recognition of handwritten cursive words". In *SPIE Symposium on Electronic Imaging Science and Technology*, San Jose, California, 1992.

- [17] J. T. Favata and G. Srikantan. "A multiple feature/resolution approach to hand-printed digit and character recognition". *International Journal of Imaging Systems and Technology*, 7:304–311, 1996.
- [18] J. Franke. "Experiments on the cenparmi data set with different structured classifiers". In *Fifth Advanced Technology Conference*, pages A167–A181, 1992.
- [19] J. Franke. "Statistical combination of multiple classifiers adapted on image parts". In *First European Meeting on Postal Technology*, pages 566–572, Nantes, France, 1993.
- [20] J. Franke and R. Lundwurm. "Different structured polynomial classifiers:performance and error effects". In *Fourth International Workshop on Frontiers in Handwriting Recognition*, pages 321–329, Taiwan, 1994.
- [21] J. Franke and M. Oberlander. "Different structured polynomial classifiers:performance and error effects". In *Second International Conference on Document Analysis and Recognition*, pages 581–584, Japan, 1993.
- [22] N. Friedman, D. Geiger, and M. Goldszmidt. "Bayesian network classifiers". *Machine Learning*, 29(2/3):131–163, 1997.
- [23] T. D. Garvey, J. D. Lowrance, and M. A. Fushler. "An interference technique for integrating knowledge from disparate sources". In *Seventh International Joint Conference on Artificial Intelligence*, pages 319–325, Vancouver, Canada, 1981.
- [24] H. Y. Hau and R. L. Kashyap. "Believe combination and propagation in a latice structured inference network". *IEEE Transactions on Systems, Man, and Cybernetics*, 20:45–57, 1990.

- [25] T. K. Ho. "Recognition of handwritten digits by combining independent learning vector quantizations". In *Second International Conference on Document Analysis and Recognition*, pages 818–821, Japan, 1993.
- [26] T. K. Ho, J. J. Hull, and S. N. Srihari. "Combination of structural classifiers". In *IAPR Workshop on Systactic and Structural Pattern Recognition*, pages 123–137, New Jersey, 1990.
- [27] T. K. Ho, J. J. Hull, and S. N. Srihari. "A regresion approach to combination of decision by multiple character recognition algorithms". In *SPIE Symposium on Electronic Imaging Science and Technology*, 1992.
- [28] T. K. Ho, J. J. Hull, and S. N. Srihari. "Desicion combination in multiple classifier systems". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(1):66–75, 1994.
- [29] Y. S. Huang, K. Liu, and C. Y. Suen. "A neural network approach for multi-classifier recognition systems". In *Fourth International Workshop on Frontiers in Handwriting Recognition*, pages 235–244, Taiwan, 1994.
- [30] Y. S. Huang and C. Y. Suen. "Combination of multiple classifiers with measurement values". In *Second International Conference on Document Analysis and Recognition*, pages 598–601, Japan, 1993.
- [31] Y. S. Huang and C. Y. Suen. "A method of combining multiple experts for the recognition of unconstrained handwritten numerals". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(1):90–94, 1995.
- [32] J. J. Hull, A. Commike, and T. K. Ho. "*Multiple algorithms for handwritten character recognition*". Concordia University, Montreal, Canada, 1990.

- [33] J. B. H. II and B. A. Pearmutter. "Evidence proof for multilayer perceptron classifiers and the bayesian discriminant function". In *1990 Connectionist Models Summer School*. Morgan Kaufmann, 1990.
- [34] A. K. Jain and R. C. Dubes. "*Algorithms for Clustering Data*". Prentice-Hall, Englewood Cliff, New Jersey, 1988.
- [35] D. W. H. Jr. and S. Lemeshow. "*Applied Logistic Regression*". John Wiley & Sons, 1989.
- [36] T. Kailath. "The divergence and bhattacharya distance measures in signal selection". *IEEE Transactions on Communication Techniques*, 15(1):52–60, 1967.
- [37] G. Kim and V. Govindaraju. "A lexicon driven approach to handwritten word recognition for real-time applications". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4):366–379, 1997.
- [38] F. Kimura and M. Shridhar. "Handwritten numerical recognition based on multiple algorithms". *Pattern Recognition*, 24(10):969–983, 1991.
- [39] J. Kittler, M. Hatef, R. P. Duin, and J. Matas. "On combining classifiers". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(3):226–239, 1998.
- [40] R. Kohavi. "Scaling up the accuracy of naive-bayes classifiers: a decision tree hibrid". In *Second International Conference on Knowledge Discovery and Data Mining*, San Mateo, California, 1996. Morgan Kaufman.
- [41] T. Kohonen. "The self-organizing map". *Proceeding of IEEE*, 78:1464–1480, 1990.
- [42] L. Lam, , and C. Y. Suen. "Increasing experts for majority vote in ocr: Theoretical considarations and strategies. In *Fourth International Workshop on Frontiers in Handwriting Recognition*, pages 245–254, Taiwan, 1994.

- [43] L. Lam, , and C. Y. Suen. "Optimal combinations of pattern classifiers". *Pattern Recognition Letters*, 16:945–954, 1995.
- [44] D. Lee and S. N. Srihari. "Handwritten digit recognition : a comparison of algorithms". In *Third International Workshop on Frontiers in Handwriting Recognition*, pages 153–162, Buffalo, New York, 1993.
- [45] D. Lee, S. N. Srihari, and R. Gaborski. "Experiments in handwritten character recognition with pattern recognition and neural network approaches". In *Fourth USPS Advanced Technology Conference*", 1990.
- [46] Y. Li and F. Yamaoka. "Integration of handwritten digit recognition results using evidential reasoning". In *Fourth International Workshop on Frontiers in Handwriting Recognition*, pages 456–463, Taiwan, 1994.
- [47] X. Ling and W. G. Rudd. "Combining options from several experts". *Applied Artificial Intelligence*, 3:439–452, 1989.
- [48] S. Madhvanath, E. Kleinberg, and V. Govindaraju. "Emperical design of multi-classifier thresholding/control strategy for recognition of handwritten street names". *International Journal of Pattern Recognition and Artificial Intelligence*, 1997.
- [49] S. Madhvanath, E. Kleinberg, and V. Govindaraju. "Holistic verification for handwritten phrases". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(12):1344–1356, 1999.
- [50] E. Mandler and J. Schurmann. "*Combining the classifiers results of independent classifiers based on the Dempster-Shafer theory of Evidence*", pages 381–393. Elseviser Science, 1988.

- [51] J. L. McClelland and D. E. Rumelhart. "*Parallel distributed processing*". MIT, 1986.
- [52] G. J. McLachlan and K. E. Basford. "*Mixture models: inference and Applications to clustering*". Marcel Dekker, Inc., 1988.
- [53] M. D. Meteish, P. Yao, and T. Stirtzinger. "A study on the use of believe functions for medical expert systems". *Journal Appl. Statistics*, 18(1):155–174, 1991.
- [54] J. R. Munkres. "*Topology, a first course*". Prentice Hall, 1975.
- [55] C. Nadal, R. Legault, and C. Y. Suen. "Complementary algorithms for the recognition of totally unconstrained handwritten numerals". In *International Conference on Pattern Recognition*, volume 1, pages 443–449, New Jersey, 1990.
- [56] T. Noumi, T. Matsui, I. Yamashita, T. Wakahara, and T. Tsutsumida. "Results of second iptp character recognition competition and studies on multi-expert handwritten numeral recognition". In *Fourth International Workshop on Frontiers in Handwriting Recognition*, pages 456–463, Taiwan, 1994.
- [57] S. J. Nowlan. "Competing experts: An experimental investigation of associative mixture models". Technical Report CRG-TR-90-5, University of Toronto, 1990.
- [58] S. J. Nowlan. "*Soft competitive adaptation: neural network by soft weight-sharing*". PhD thesis, 1991.
- [59] T. W. O’Gorman and R. F. Woolson. "Variable selection to discriminate between two groups: Stepwise logistic regression or stepwise discriminant analysis". *The American Statistician*, 45:187–193, 1991.
- [60] Y. Pearl. "*Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Interference*". Morgan Kaufmann Publishers, San Mateo, California, 1988.

- [61] B. Plessis, A. Sicsu, L. Hautte, and etc. "A multi-classifier combination strategy for the recognition of handwritten cursive words". In *Second International Conference on Document Analysis and Recognition*, pages 642–645, Japan, 1993.
- [62] S. J. Press and S. Wilson. "Choosing between logistic regresion and discriminant analysis". *Journal of American Statistical Association*, 73:699–705, 1978.
- [63] M. Ramoni and P. Sebastiani. "*Bayesian Methods for Intelligent Data Analysis*". Springer, New York, 1999.
- [64] M. D. Richard and R. P. Lippmann. "Neural network classifiers estimate bayesian a postiriori probabilities". *Neural Computation*, 3:461–483, 1978.
- [65] G. Rogaua. "Combining the results of several neural network classifiers". *Neural Networks*, 7(5):777–781, 1994.
- [66] W. Rudin. "*Real and complex analysis*". McGraw-Hill Companies, 1986.
- [67] M. Sabourin, A. Mitiche, D. Thomas, and G. Nagy. "Classifier combination for handprinted digit recognition". In *Second International Conference on Document Analysis and Recognition*, pages 163–166, Japan, 1993.
- [68] G. Shafer. "*A mathematical theory of evidence*". Princeton, 1976.
- [69] S. N. Srihari. "Recognition of handwritten and machine printed text for postal address interpretation". *Pattern Recognition Letters*, 14:291–302, 1993.
- [70] G. Srikantan. "*Image sampling rate and image pattern recognition*". PhD thesis, SUNY at Buffalo, 1994.
- [71] U. Srinivasan. "Polynomial discriminant method for handwritten digital recognition". Technical report, SUNY at Buffalo, Sep 1989.

- [72] C. Y. Suen, R. Legault, C. Nodel, M. Cheriet, and L. Lam. "Building a new generation of handwriting recognition systems". *Patterns Recognition Letters*, 16(1):66–75, 1994.
- [73] C. Y. Suen, C. Nodel, R. Legault, T. Mai, and L. Lam. "Recognition of totally unconstrained handwritten numerals based on the concept of multiple experts". In *First International Workshop on Frontiers in Handwriting Recognition*, pages 131–143, Montreal, Canada, 1990.
- [74] C. Y. Suen, C. Nodel, R. Legault, T. Mai, and L. Lam. "Computer recognition of unconstrained handwritten numerals". *Proceedings of IEEE*, 80(7):1162–1180, 1992.
- [75] T. Suzuki, H. Nishida, Y. Nakajima, H. Yamagama, M. Tachikawa, and G. Sato. "A fast and high-performance system for handwritten character recognition by efficient combination of multiple classifiers". In *International Association for Pattern Recognition Workshop on Document Analysis Systems*, pages 175–189, Kaiserslautern, Germany, 1994.
- [76] H. Tahani and J. M. Keller. "Information fusion in computer vision using the fuzzy integral". *IEEE Transactions on System, Man, and Cybernetics*, 20(3):733–741, 1990.
- [77] H. Takahashi and T. D. Griffin. "Recognition enhancement by linear tournament verification". In *Second International Conference on Document Analysis and Recognition*, pages 585–588, Japan, 1993.
- [78] J. T. Tou. *Feature selection for pattern recognition systems*. Academic Press, 1969.

- [79] R. Tubbs and W. O. Alltop. "Measures of confidence associated with combining classification results". *IEEE Transactions on System, Man, and Cybernetics*, 21:690–692, 1990.
- [80] E. Wan. "Neural network classification: a bayesian interpretation". *IEEE Transactions on Neural Networks*, 1(4), 1990.
- [81] D. H. Wolpert. "Stacked generalization". *Neural Networks*, 5(2):241–260, 1992.
- [82] L. Xu and M. I. Jordan. "Em learning on a generalized finite mixture model for combining multiple classifiers". In *World Congress on Neural Networks*, volume 4, pages 227–230, 1993.
- [83] L. Xu, A. Krzyzak, and C. Y. Suen. "Methods of combining multiple classifiers and their application to the handwritten numerical recognition". *IEEE Transactions on System, Man, and Cybernetics*, 22(3):418–435, 1992.
- [84] F. Yamaoka, Y. Lu, A. Shaout, and M. Shridhar. "Fuzzy integration of classification results in a handwritten digit recognition system". In *Fourth International Workshop on Frontiers in Handwriting Recognition*, pages 255–264, Taiwan, 1994.
- [85] L. A. Zadeh. "A simple view of dempster-shafer theory of evidence and its implication for the rule of combination". *Artificial Intelligence Magazine*, pages 85–90, 1986.