

RECURSIVE HANDWRITING RECOGNITION USING EPISTEMIC UTILITY THEORY

Jaehwa Park and Venu Govindaraju

Center of Excellence for Document Analysis and Recognition
Department of Computer Science and Engineering
State University of New York at Buffalo
Amherst, NY 14260, USA
{jehwap,govind}@cedar.buffalo.edu

ABSTRACT

In most handwriting recognition systems, the processing steps of segmentation, recognition, decision making and postprocessing are serially taken. These systems usually use the resources exhaustively in each stage of the serial engine. Every stage is tuned to maximize global performance.

However, such a paradigm is prone to miss the global optimal point because the individual stages often do not have all the information available. A more interactive model is desirable. In this paper we describe such a model extending the concepts of satisficing decision mechanisms. Our experiments show proof of validity of the proposed model.

1. INTRODUCTION

General processing steps of handwriting document recognition are *preprocessing segmentation, recognition, decision making* and *postprocessing*. *Preprocessing* is primarily related to image processing operations such as normalization to remove irregularities of handwriting. *Segmentation* separates an image into meaningful units of recognition. *Recognition* is the application of classification algorithms. Independent contextual information is used to reorganize and enhance as a *postprocessing* step.

Most architectures of handwriting recognition systems are monotonically linear. Starting from a scanned image to recognition every step is sequential. Further each step is usually independent. Resources are exhaustively used at each stage of the serial engine without consideration of environmental changes. The idea is to maximize the global performance at each stage.

Recognition engines operate in a static model where a uniform classification algorithm is applied on all classes with a fixed feature space that is optimally

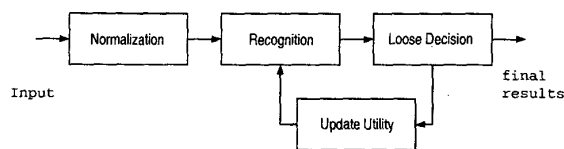


Figure 1: recursive recognition

selected from a training set. Decision making is designed independently to provide one and only one decision output which is represented as the true answer. Even if multiple recognition experts are used, it is not interactive nor dynamic. [1] However, such a paradigm is prone to miss the optimal point because individual stages usually do not have all the information.

Cognitively speaking, humans do not follow such monotonic paradigms. They can alter course of action dynamically and select the utilities relative to their goals.

In order to extend traditional recognition systems to dynamic and recursive systems, two concepts must be introduced: (i) altering action using feedback and environment and (ii) evaluating and regulating terminal conditions actively.

In this paper we introduce dynamic and recursive recognition paradigm with: (i) a multiresolution hierarchical optical character recognizer, [2] (ii) *satisficing decision maker* (epistemic utility theory [3]) and (iii) lexicon complexity.

2. RECURSIVE RECOGNITION

If recognition results are made available by a classifier as evidence, the goal of decision making is to meet the performance requirements that are imposed by the users of the recognition system.

Most lexicon driven handwriting recognition sys-

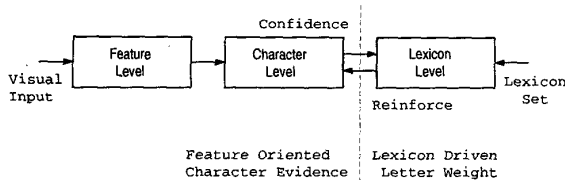


Figure 2: Lexicon-Adaptation Recognition

tems use an OCR to compute the distance to test character images from prototypes that can match. Lexicons are used to confine the possible set of prototypes. (Figure 2).

During the recognition stage, a classifier uses all resources in order to make the most informed decision. Cost-optimal results can be achieved by interaction between recognition and decision making.

Epistemic utility theory [3] provides a means of quantifying the value of avoiding error and the value of acquiring new information. The first principle, *epistemic utility*, regulates evaluation of information regardless of its truth-value, and the second principle, *credal probability*, regulates the credibility of information regardless of its worth. These two principles lead to of dynamic decision making.

Epistemic utility theory is commonly used in the theory of *satisficing decision* [4]. The notion of truth and information is re-interpreted in a practical setting. The notions of credal probability and epistemic utility are extended to *accuracy* (conformity to a given standard) and *liability* (susceptibility or exposure to something undesirable).

We use the basic *satisficing recognition* concept in handwriting recognition by using the confidence measures returned by OCRs as the metric

2.1. Utility Theory

Let Ψ denote the full truth set and H a pre-selected subset of Ψ including the null set ϕ . H consists of possible true hypotheses.

$$H = \{\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_{K-1}\} \quad (1)$$

In handwriting recognition, Ψ represents all possible combinations of characters and H is a meaningful lexicon set of K elements. H is narrowed from Ψ independently. If we assume that all lexicons are unique and are not subsets of other lexicons then \mathbf{h} represents a class array of character sets.

$$\mathbf{h} = [h_0, h_1, \dots, h_{L-1}] \quad (2)$$

where the class set is an extended class including the null set ϕ ,

$$C = \{c_0, c_1, \dots, c_{N-1}, \phi\} \quad (3)$$

where $h \subset C$

The epistemic utility function of each component is the convex combination of accuracy and liability [4].

$$\tilde{u}(\mathbf{h}/H) = aQ(\mathbf{h}/H) + (1-a)(1-M(\mathbf{h}/H)) \quad (4)$$

where $a \in [1, 0]$. After positive linear transformation the function is

$$u(\mathbf{h}/H) = Q(\mathbf{h}/H) - bM(\mathbf{h}/H) \quad (5)$$

where

$$u(\mathbf{h}/H) = [u_0, u_1, \dots, u_{L-1}] \quad (6)$$

$$u_i = u(h_i/H) \quad (7)$$

The parameter b is termed as the *index of rejectivity*. Large values of b imply an increased willingness to reject hypotheses. $b \cong 0$ implies minimal concern for liability, $b \cong 1$ implies equal concern for accuracy and liability. $Q(\mathbf{h}/H)$ is a distance measure representing the accuracy support associated with \mathbf{h} at given H , and $M(\mathbf{h}/H)$ characterize the liability exposure of \mathbf{h} in H .

If multiple symbols are assigned to an h , for instance, allowing folding of lower and upper characters, each epistemic utility value is defined by

$$u_i = \text{MAX}(u(c))|_{vc \in h_i} \quad (8)$$

If we set the overall epistemic utility value, \bar{h} of a hypothesis \mathbf{h} as the average value of each epistemic utility

$$\bar{h}(\mathbf{h}/H) = \frac{1}{L} \sum_{i=0}^{L-1} u_i \quad (9)$$

In a segmentation oriented lexicon driven word recognizer [5], the distance of each test pattern from the prototype can be used as the value of accuracy support $Q(h_{ij})$ of hypotheses h_{ij} where i is lexicon index and j is index of symbol sequence. In the above equation, if we set *index of rejectivity* b as zero, \bar{h} becomes the method of finding the average distance.

The liability, $M(h_{ij})$, represents the information value of rejecting h_{ij} . In handwriting recognition, lexicon complexity measures the importance of each character within a lexicon set. If every lexicon entry starts with same character, the first character is less informative. Inter-class confusions reflect imperfections of the OCR. Inter-character distance is the separation characteristic of an OCR in terms of regulating confidence values.

2.2. Revision of Utility

The essence of decision-making strategy lies in rejecting only the "bad" options and leaving other options for further consideration. Priority is in avoiding error rather than minimizing the probability of error. If we recursively reject only bad options, we will eventually reach the satisfying conditions of termination.

There are two ways in which a recognizer may wish to modify its decision boundary without adopting new process schedule (Eq 5). First, renew the accuracy support of hypotheses using a more accurate OCR. Second, update the liability exposure by discarding hypotheses which are no longer believed to be true.

Since accuracy support basically depends on the distance measures returned by the OCR, it is renewable by applying more accurate classification. Liability exposure is also renewable by considering complexity of lexicon set and confusions of OCR.

If the decision maker discards some hypotheses from lexicon set as false, the matching complexity is lowered and the liability exposure values are updated. If the decision maker no longer updates the liability exposure, evidence should be upgraded to expand the decision boundary.

In this control scheme, we can rewrite the epistemic utility function in Eq 5 as a parameter dependent recursive equation,

$$u(\mathbf{h}/H(t)) = Q(\mathbf{h}/H(t-1), \eta) - bM(\mathbf{h}/H(t-1), \eta) \quad (10)$$

where η is a control variable of OCR and $H(t)$ denotes the narrowed hypotheses set at time t by the decision maker.

3. LEXICON COMPLEXITY

The object in most word recognizers is to find the best matching between a symbol string and image segmented strokes. In [5], a two step dynamic matching has been implemented. Comparisons between feature vectors extracted from possible combinations of segments and reference prototypes of codewords are made to find the best match. In the second phase, a global optimum path is obtained by dynamic programming using the scores from the previous step.

However, the successful recognition cannot always be realized due to limitations of information and resources during the training. This can be overcome by knowledge of environmental facts. Entropy of a lexicon (lexicon complexity) and the reliability of an OCR module serve to capture the dynamism of the environment.

Lexicon Set 1: West Central Street
East Central Street
Sunset Avenue

Lexicon Set 2: West Central Street
West Main Street
Sunset Avenue

Lexicon Set 3: West Central Street
West Central Avenue
Sunset Avenue

Figure 3: lexicon complexity

Let us consider an input image and three different lexicon sets (Figure 3). Initially, all three lexical entries and their constituent characters are of equal importance. But after rejecting 'Sunset Avenue' in set 1 the first two characters are more informative than the other characters. In lexicon set 2, the substring 'Central' is most informative than the entire phrase to choose truth as 'West Central Street'. In lexicon set 3, the suffix 'Street' is most informative.

If the OCR do not provide strong separation between the substrings of 'We' and 'Ea', it is not confident in accepting 'West Central Street' as the final decision. If the OCR provides a strong classification between 'Street' and 'Avenue', 'West Central Street' can be accepted as final result even if the absolute confidence score is same. To maximize success, measuring lexicon complexity in terms of character importance must be considered.

The distance between two finite character arrays can be measured by *Edit Distance* [6]. Given two strings edit distance is defined as the minimum cost of transforming one string to the other through a sequence of weighted edit operations. The edit operations are usually defined in terms of *insertion*, *deletion*, and *substitution* of symbols.

3.1. Edit Cost

Consider two different lexicon hypotheses of \mathbf{h}_i and \mathbf{h}_j over a finite alphabet length L_i and L_j each as,

$$\mathbf{h}_i = [h_{i,0}, h_{i,1}, \dots, h_{i,L_i-1}] \quad (11)$$

$$\mathbf{h}_j = [h_{j,0}, h_{j,1}, \dots, h_{j,L_j-1}] \quad (12)$$

Let us assume that a stroke segmentation algorithm is applied to a given image and a stroke segment array \mathbf{s} of length R is generated as

$$\mathbf{s} = [s_0, s_1, \dots, s_{R-1}] \quad (13)$$

Also assume that a dynamic programming based word recognizer (similar to [5]), finds the best matching path

within a segment array conforming to rules such as (i) at least one character matches one stroke segment, (ii) one character can only match within a window, and (iii) lexical entries which are smaller than R or the maximum allowable number of characters of stroke array are discarded.

We can define the edit operations without altering the framework of matching stroke segments as follows

- *nullifying*: $t(h_{i,m} \rightarrow \phi/s_{a\sim b})$
- *substitution*: $t(h_{i,m} \rightarrow h_{j,n}/s_{a\sim b})$
- *splitting*: $t(h_{i,m} \rightarrow [h_{j,p}, \dots, h_{j,q}]/s_{a\sim b})$

Nullifying is the operation matching a character to certain boundary of stroke segments. *Substitution* is the replacement of a character matching to another with exact same matching bound, and *splitting* is matching of a character into several characters.

The overall matching distance of a lexical entry primarily depends on the distance measurement between extracted feature vector of one or more combined stroke segments and corresponding prototypes. Let $\varepsilon(x \rightarrow y/s_{a\sim b})$ be the elementary transform cost function from x to y and $\delta(x/s_{a\sim b})$ be the OCR distance function of x within stroke segments s_a to s_b , where where $x, y \in C$ and $0 \leq a \leq b \leq R$.

The transform cost functions are defined as

$$\varepsilon(h_{i,m} \rightarrow \phi/s_{a\sim b}) = -\delta(h_{i,m}/s_{a\sim b}) \quad (14)$$

$$\varepsilon(h_{i,m} \rightarrow h_{j,n}/s_{a\sim b}) = \delta(h_{j,n}/s_{a\sim b}) - \delta(h_{i,m}/s_{a\sim b}) \quad (15)$$

$$\varepsilon(h_{i,m} \rightarrow [h_{j,p}, \dots, h_{j,q}]/s_{a\sim b}) = \text{MIN}(\delta(h_{j,p}/s_{a\sim c_1}) + \delta(h_{j,p+1}/s_{c_1+1\sim c_2} + \dots + \delta(h_{j,p}/s_{c_{q-p-1}+1\sim b}) - \delta(h_{i,m}/s_{a\sim b})) | \forall a \leq c_1 \leq \dots \leq c_{q-p-1} \leq b \quad (16)$$

where n, p, q are indexes of eligible target lexicon array.

3.2. Matching Distance

Given a transforming path from \mathbf{h}_i to \mathbf{h}_j , the matching distance has been defined directly in terms of paths (or traces)

$$t(\mathbf{h}_i \rightarrow \mathbf{h}_j) = \mathbf{t}_{i,j} = [t_{i,j,0}, t_{i,j,1}, \dots, t_{i,j,L_i-1}] \quad (17)$$

where t is a transforming operation of symbols within the framework of \mathbf{h}_i defined as before.

The associated transform cost to paths are

$$w(\mathbf{t}_{i,j}) = \sum_{k=0}^{k=L_i-1} \varepsilon(t_{i,j,k}) \quad (18)$$

Let $\bar{\mathbf{t}}_{i,j} = [\bar{t}_{i,j,0}, \bar{t}_{i,j,1}, \dots, \bar{t}_{i,j,L_i-1}]$ be the transform path of minimum cost from \mathbf{h}_i to \mathbf{h}_j . Then the matching distance from \mathbf{h}_i and \mathbf{h}_j is defined as

$$d(\mathbf{h}_i \rightarrow \mathbf{h}_j) = \text{MIN}(w(\mathbf{t}_{i,j})) \quad (19)$$

$$= w(\bar{\mathbf{t}}_{i,j}) \quad (20)$$

The overall transform cost of a character hypothesis $h_{i,*}$ is defined by the sum of all transform cost to all other lexicons. The overall transform cost of $h_{i,*}$ in a lexicon set denoted as $\rho(h_{i,*})$ becomes

$$\rho(h_{i,*}) = \sum_{j=0, j \neq i}^{j=K-1} \varepsilon(\bar{\mathbf{t}}_{i,j,*}) \quad (21)$$

The liability exposure of a lexicon hypothesis $M(\mathbf{h}_i)$ can be evaluated by the average of sum of each transform cost to all other lexicon hypotheses as

$$M(\mathbf{h}_i) = \frac{1}{K-1} \rho_{i,j}(\mathbf{h}_i) \quad (22)$$

The average OCR distance along the lexical entry is the measurement of the best path within a stroke segment array. This transform distance of paths reflects the entropy of each hypothesis projected into OCR feature space within a lexicon set.

4. EXPERIMENT

The recursive recognition control strategy described is applied to dynamic programming based lexicon driven word recognizer described in [5]. The HOCR [2] is used for experiments. The basic performance of the word recognizer with maximum resource usage is 91% (top choice).

Figure 4 shows the dynamic characteristic of recursive recognition when only the usage of resources change without considering the liability between lexicons. As more resources are used, the average minimum distances grows. The separability between lexical entries also larger. If the usage of sub-nodes are increased continuously the gap between classes start to decrease at a certain point (because of dimensionality problem). However the behavior of gaps between lexical entries are unstable if the recognizer fail to separate strongly (Figure 4-(b)). The processing time is linearly proportional to the number of sub-images used.

Figure 5 shows the dynamic characteristic of adjusted distance by utility between lexical entries. The index of rejection b is 0.7, the used nodes are 96 and lexicon size is 100. If the rejected lexical entries become large, the gap between the first and second choices becomes large.

5. CONCLUSION

The basic *satisficing recognition* concept is developed by adapting epistemic utility theory. The accuracy support and liability exposure functions are formulated in terms of Transform Distance and class confusion of OCR. Also these two functions are modified into a recursive form which depends on the decision-maker and the separability threshold of OCR. The experiments show the proof of validity of the proposed model.

6. REFERENCES

- [1] Y.S. Huang and C.Y. Suen. A Method of Combining Multiple Experts for the Recognition of Unconstrained Handwritten Numerals. *IEEE Trans. Pattern Recognition and Artificial Intelligence*, 17(1):90-94, 1995.
- [2] J. Park V. Govindaraju and S.N. Srihari. OCR in A Hierarchical Feature Space. In *IEEE International Conference on Systems, man, and Cybernetics, San Diego, CA*, pages 4324-4329, Oct. 1998.
- [3] I. Levi. *Decisions and Revisions*. Cambridge Univ. Press, London, U.K., 1984.
- [4] M.A. Goodrich W.C. Stirling and R.L. Frost. A Theory of Satisficing Decisions and Control. *IEEE Trans. System, Man and Cybernetics - Part A*, 28(6):763-779, 1998.
- [5] G. Kim and V. Govindaraju. A Lexicon Driven Approach to Handwritten Word Recognition for Real-Time Applications. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 19(4):366-379, 1997.
- [6] R.A. Wagner and M.J. Fischer. The String-to-String Correction Problem. *J. Assoc. Comput. Mach.*, 21(1):168-173, 1974.

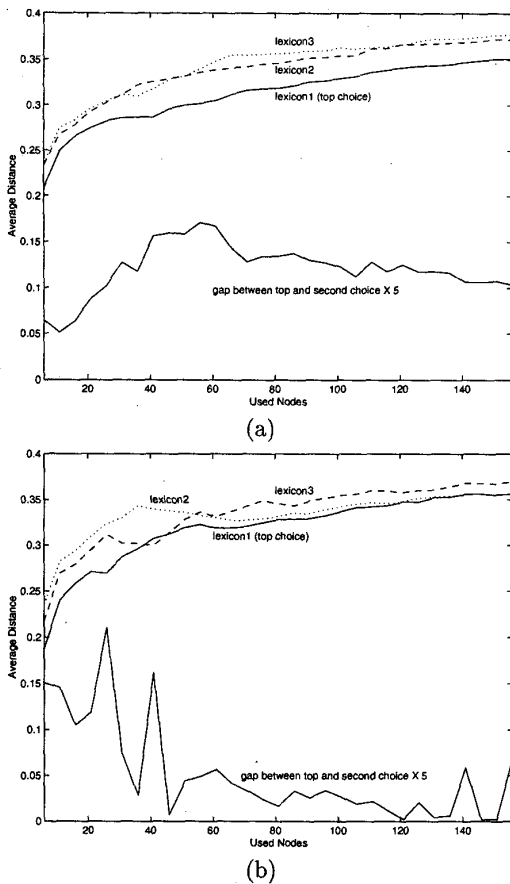


Figure 4: Accuracy Support Changes - (a) success case
(b) failure case

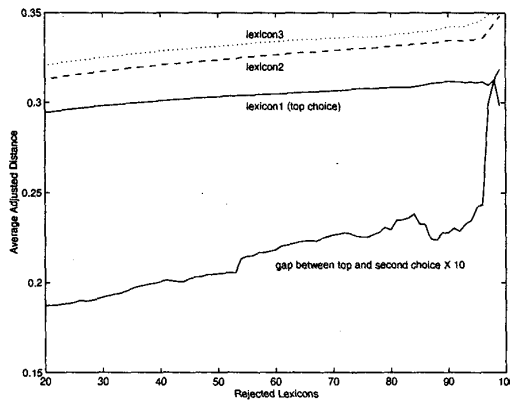


Figure 5: Adjusted Distance Changes