

On the Structure of Hidden Markov Models

K. T. Abou-Moustafa^{a,b}, M. Cheriet^b and C. Y. Suen^a

^a *CENPARMI, Concordia Univ., GM-606, 1455 de Maisonneuve W., Montréal, H3G 1M8, QC, Canada*

^b *LIVIA, École de Technologie Supérieure, Univ. de Québec, 1100 Notre-Dame W., Montréal, H3C 1K3, QC, Canada*

Abstract

This paper investigates the effect of HMM structure on the performance of HMM-based classifiers. The investigation is based on the framework of graphical models, the diffusion of credits of HMMs and empirical experiments. Although some researchers have focused on determining the number of states, this study shows that the topology has a stronger influence on increasing the performance of HMM-based classifiers than the number of states.

Key words: HMM structure, graphical models, credits diffusion, Ocham's Razor, K-Means clustering

1 Introduction

Hidden Markov Models (HMMs) (Baum et al, 1966) are a class of stochastic processes that is capable of modeling time-series data. They belong to a larger class of models known as generative models. Though generative models are tools for data modeling, in the literature, HMMs were used in many classification problems such as speech recognition (Rabiner, 1989; Baker, 1975), handwritten word recognition (Yacoubi et al, 1999), object recognition (Cai et al, 2001), gesture recognition (Kim et al, 2001), bioinformatics (Baldi et al, 1998) and modeling of biological sequences (Karplus et al, 1997).

Designing an HMM for data modeling to be a part of an HMM-based classifier means determining the structure (the number of states, and the topology) of the model. The topology in this context is meant to be the connections (transitions) between the states. The structure affects the modeling capability considerably and consequently the performance of the classifier. An estimation of the weight of each factor on the performance can point out to the main factor

affecting the performance and consequently can lead to an improvement in the selection of values of this factor. Although some researchers focused on the problem of number of states and the topology, the goal of this paper is to investigate the effect of the number of states and the topology, each separately, on the performance of HMM-based classifiers. Our research shows that the topology can improve the modeling capability greatly.

The investigation is based on 1) linking the theoretical results from model selection for graphical models with the diffusion of credits in Markovian models and, 2) supporting the results with empirical experiments for the recognition of unconstrained handwritten digits. The experiments compared the performance obtained from classifiers with different structures.

It is worth calling the readers' attention to two issues regarding this work. First, though HMMs are usually trained with time-series data with a long signal duration, such as the applications mentioned above, the experiments used isolated handwritten digits that have a short signal duration. The reason for that is to treat the HMM-based classifier like other classifiers such as Multi Layer Perceptrons (MLPs) and Support Vector Machines (SVMs) without any constraints on the data. Second, the goal of the paper is not to introduce a new state-of-the-art recognition result on the MNIST database using HMMs, but rather, to compare the performance of HMM-based classifiers under different structure conditions. It is well known from the literature that classifiers such as SVMs and MLPs can achieve very high recognition results (Dong, 2003; Simard et al., 2003; Liu et al., 2003) on this database.

For complete references on HMMs, the reader is required to read (Rabiner, 1989; Bengio, 1999). The paper uses the basic compact notation of HMMs defined as $\lambda = (A, B, \pi)$ where λ is the hidden Markov model, A is the transition probability matrix, B is the observation probability matrix and π is the initial state probability

The rest of the paper is organized as follows: Sections 2 and 3 review related work in the literature and HMMs respectively. Section 4 discusses the effect of the structure on the modeling capability of HMMs. Section 5 describes the experiments and results and finally section 6 concludes the paper.

2 Related Work

A work directly related to this investigation is the problem of optimizing HMM structure in two forms; 1) application dependent methods, and 2) application independent methods. Application dependent methods use a priori knowledge from the application domain such as (Yacoubi et al, 1999), where they used

information extracted from a character segmentation process to build a special HMM structure (number of states and the topology), (De Britto et al, 2001) modified the left-to-right model to enhance the performance of his proposed framework for numeral strings and (Lee et al, 2001) fixed the topology to be a left-to-right one and determined the number of states by reflecting the structure of a target pattern. The major drawback of these methods is that they are designed for specific applications and can not be generalized to others.

On the other hand, application independent methods, although are more promising, yet they are not popularized. These methods include the work of (Stolcke et al, 1992) and (Brants, 1996) where they proposed an incremental learning for the structure based on state merging and splitting, i.e. the structure is changed as new evidence is added to the model; (Lien, 1998) proposed a general method to determine the number of states and the connections between states in discrete left-to-right HMMs. Recently, Bicego *et al* focused only on determining the number of states using probabilistic bisimulation (Bicego et al, 2001) and sequential pruning using Bayesian Information Criterion (BIC) (Bicego et al, 2003). Model selection approaches were also investigated for this purpose and recently (Biem, 2003) proposed a Discriminative Information Criterion (DIC) framework and used it to optimize the HMM structure.

Different approaches for structure optimization can also be found. (Lyngso et al, 1999) focused on comparing HMMs in terms of state emission probabilities, (Bahlman et al, 2001) used Bayesian estimates of HMM states as a criterion for selecting HMMs and (Balasubramanian, 1993) selected HMMs based on equal probabilities of observation sequences only. By examining the above literature, and except for the work Stolcke and Brants, most of these methods use the left-to-right topology and the optimization targets only the number of states and the number of mixtures in cases of continuous HMMs. However, according to this investigation, the number of states may not affect the performance after a certain limit, but it can reduce the computational time for training and testing, while the topology can considerably affect the performance of HMMs.

3 HMM Structure

In many applications that use HMMs, the number of states is manually pre-determined prior to training. The connections between states, (topology) is determined by setting non-zero probabilities in the A matrix prior training. During training, the EM (Baum-Welch) algorithm improves the estimates of these probabilities from the data. Note that the EM algorithm can not set 0 or 1 (can approach 0 or 1) probabilities in the A matrix, therefore it can not be seen as an algorithm that learns the topology. In the following, we inves-

tigate the effect of the topology on the performance of HMM-based classifiers through two different perspectives: (1) using the graphical models framework and, (2) using the diffusion of credits while learning Markovian models.

3.1 Bayesian formulation for model selection

Determining the number of states and the topology of HMMs can be viewed as a model selection problem. The problem can be formulated as follows. Given the training set of examples Ψ and a criterion function Υ for the quality of the model on the data set Ψ , choose a model from a certain set of models, in such a way to maximize the expected value of this criterion function on new data (assumed to be sampled from the same unknown distribution from which the training data was sampled) (Bengio, 1999).

HMMs can be viewed as a special case of Graphical Models (Heckerman, 1996; Murphy, 2001). Model selection is one of the main problems in graphical models and much work has been introduced regarding this problem. The Bayesian approach, one of the main approaches for model selection, is a fundamental approach for model selection in graphical models. Following this approach means encoding the uncertainty about the structure of the HMM by using a discrete variable whose states correspond to the possible HMM-structure hypotheses S^h and assessing it the a priori density $P(S^h)$. Given the training example set Ψ for the model λ and augmenting the model parameters A, B, π in a single parameter vector θ , the problem would be computing the posterior distribution for the HMM structures. This can be formulated as follows using Bayes theorem:

$$P(S^h|\Psi) = \frac{P(\Psi|S^h)P(S^h)}{P(\Psi)} \quad (1)$$

where $P(\Psi)$ is a constant that does not depend on the network structure.

The maximum likelihood structure would be the complete graph (Murphy, 2001), i.e. the full ergodic model, since this has the greatest number of parameters, and hence can achieve the highest likelihood. On the other hand this increases the model's complexity and will let the model overfit the training data resulting in a poor generalization. In fact, the marginal likelihood in 1 plays an important role to prevent this overfit. From the definition of the marginal likelihood:

$$P(\Psi|S^h) = \int P(\Psi|S^h, \theta)P(\theta|S^h) d\theta \quad (2)$$

it automatically penalizes more complex structures since they have more pa-

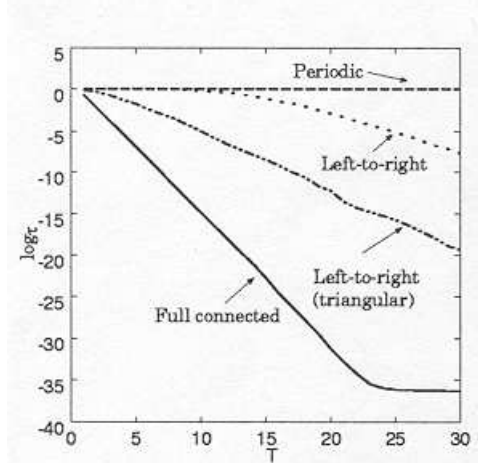


Fig. 1. Convergence of Dobrushin's coefficient for 4 different topologies.

rameters and hence cannot give as much probability mass to the region of space where the data actually lies. In other words, a complex model is less believable and hence less likely to be selected. This phenomenon is known as Ockham's Razor (Murphy, 2001) which favors simple models over complex ones. It can be seen that though the number of states may be fixed, the topology can affect the modeling capability in a serious way.

3.2 Diffusion of credits in Markovian models

The work in (Bengio et al, 1995) investigated the problem of diffusion in homogeneous and non-homogeneous HMMs and its effect on learning long term dependencies. Training HMMs requires propagating forward and backward probabilities and taking products of the transition matrix. Therefore, two types of diffusion exist, diffusion of influence in the forward path and diffusion of credit in the backward phase of training. The paper (Bengio et al, 1995) studied under which conditions these products of matrices will converge to a lower rank, thus harming learning long term dependencies. The difficulty of learning was measured by using the Dobrushin's ergodicity coefficient (Senta, 1986) defined as follows:

$$\tau(A) = \frac{1}{2} \sup_{i,j} \sum_k |a_{ik} - a_{jk}| \quad (3)$$

where $A = \{a_{ij}\}$ is the transition probability matrix. It was shown that in all cases, while training HMMs, the ergodicity coefficient will converge to 0 indicating a greater difficulty in learning, but the rate of convergence depends on the topology. Figure 1 (Bengio et al, 1995) shows the convergence of 4 HMMs with the same number of states but with different topologies. It can

be seen that the full ergodic model has the fastest convergence rate and that simpler models are slower. The final conclusion is that in order to avoid any kind of diffusion, most transition probabilities should be deterministic (0 or 1 probability). The result coincides with the Ockham's Razor result obtained from the previous section and both prefer simple topologies over complicated ones.

4 Experiments

We were interested in investigating experimentally how the number of states and the topology can affect the performance of an HMM-based classifier. Two types of experiments were carried out, one to study the effect of number of states on the performance, and the other to study the effect of the topology on the performance.

4.1 *The dataset and feature extraction*

The dataset used in the experiments consists of images of unconstrained handwritten digits from the MNIST database (LeCun, 1998) which has a training set of 60,000 samples and a test set of 10,000 samples from approximately 250 writers. The digits are cropped and scaled to be contained in a 20x20 pixels images. The gray level values of the images were normalized to be from 0 to 1. The time series data were extracted from the digits using the sliding window technique (Cornell, 1996) with a width of 3 pixels, height equals the image height and an overlap of 2 pixels. A feature vector is extracted from each window by computing the average gray level value in each row of the window, i.e. the sum of gray level pixels in each row divided by the window width. This resulted in an observation sequence length of 18 vectors from each image.

4.2 *HMM density type, initialization and codebook size*

The experiments were conducted using discrete HMM (DHMM) -based classifiers where each consisted of 10 DHMMs. The number of states for each model was determined according to the goal of the experiment. Two topologies were used in the experiments, the left-to-right with self-state transitions (no jumps), and the ergodic topology. For the code book, the vector quantization (Gray, 1984) algorithm was used to construct seven different code books (16, 32, . . . , 1024). The initial parameters for B in all experiments were set using a uniform distribution. In our original investigation, all the experiments were conducted

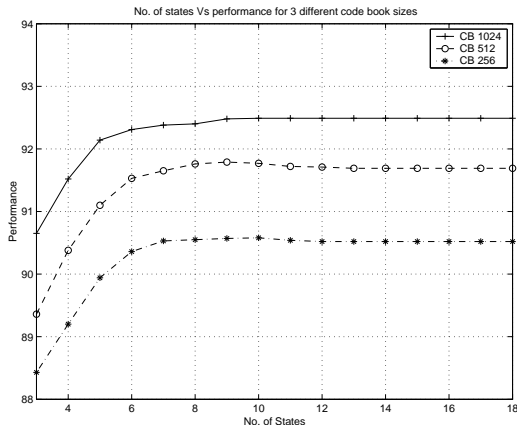


Fig. 2. The relation between performance and the number of states with different code book sizes.

using the seven code books and with several initializations for the A matrix as will be shown later. However, due to the following reasons: 1) space limitation, 2) avoid redundancy, and 3) similarity of results and conclusions, we selected the clearest of these experiments for illustration.

4.3 Studying the effect of number of states

In studying the effect of number of states, two experiments were conducted. The first experiment used HMMs with a left-to-right topology and all models had an equal number of states. The experiment studied the relation between the performance and the increase in the number of states in the classifier. The second experiment studied the performance of classifiers with a varying number of states in each model. It compared the performance between models with an equal number of states and models with a varying number of states.

4.3.1 Experiment 1

This experiment was conducted using the seven code books, and for each experiment, the A matrix was initialized using 3 different initializations; (0.5 & 0.5, 0.7 & 0.3 and 0.9 & 0.1) for the ij and ii transitions respectively. Figure 2 illustrates the results for *Experiment 1* using three code books and the first initialization for the A matrix. It can be seen that increasing the number of states can increase the performance up to a certain limit, after that a saturation is reached whenever more unnecessary states are added. However, the saturation may be accompanied by a slight drop in the performance.

The saturation may be explained as follows. The number of states N , is the number of values that the hidden variable can take and accordingly the emis-

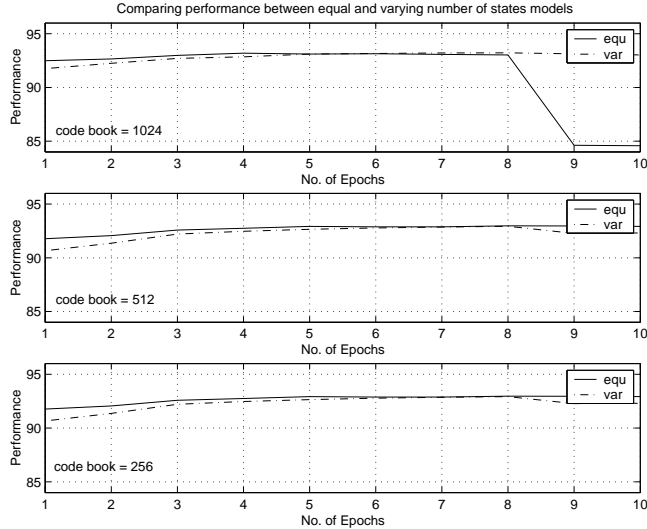


Fig. 3. Performance comparison between models with equal (EQU) and varying (VAR) number of states with different codebook sizes.

sion of symbols change. Let the true (unknown) number of values of the hidden variable be N_0 . If $N \ll N_0$ poor modeling will result and hence a classifier with poor performance. If $N \gg N_0$, additional states will introduce redundancy with no effect on the modeling capability and hence the performance is saturated. Adding more unnecessary states increases the complexity (time and computation) with no effect on the performance.

4.3.2 Experiment 2

The goal of the experiment was to measure the performance of classifiers with a different number of states in each model to see how comparable they are with classifiers having all models with an equal number of states. Two HMM classifiers were used. According to the previous experiments, the first classifier had 10 states per model, the second classifier had a different number of states in each model. Determining the number of states in each model will be described in the next subsection. As *Experiment 1*, this experiment was also conducted using the seven code books and the three different initializations. Figure 3 illustrates the results of this experiment for three code book sizes and the first initialization. Models with an equal number of states are referred as (EQU) and models with a varying number of states are referred as (VAR).

Figure 3 shows clearly how models with a varying number of states can achieve almost the same performance of models with an equal number of states with the advantage of a smaller number of states but paying the price of more epochs. The total number of states in the EQU models is 100, and the total for VAR models is 70 states. Achieving the same performance with a smaller number of states means a considerable reduction in complexity when it comes

Table 1
The number of states of each model.

<i>Model</i>	0	1	2	3	4	5	6	7	8	9
<i>No. of States (3-5)</i>	5	5	5	5	4	4	4	5	3	4
<i>No. of States (3-9)</i>	6	5	8	8	9	6	8	8	3	9

to large classification problems. However, as followed in the literature (Yacoubi et al, 1999; Augustin et al, 1998), a guaranteed performance with an easy design would be an HMM classifier with an equal number of states for all models. In Figure 3, it is worth mentioning that the drop seen in the first graph is experienced in the other graphs for the EQU and VAR models but in late epochs not shown in the graphs. The reason for the drop is due to the overfit of models on the data and due to the diffusion of credits while learning.

4.3.3 Determining the number of states

As mentioned earlier, the number of states is usually fixed (manually predetermined). Exceptions are models that use automatic clustering algorithms that determine the number of states and their outputs, but this still leaves out the topology (Brants, 1996; Theodoridis et al, 1999). Clustering sequential data while neglecting the variations of the time factor, tends to discover the underlying structure of the data given that the number of clusters is known. To determine the number of states using clustering, we proposed the use a cluster validity index (Bezdek et al, 1998) to measure the goodness of different clustering configurations and then select the best number of clusters according to this cluster validity index.

In the experiments, the K-Means algorithm (Duda et al, 2001) was used to cluster the sequential data of each model. The algorithm was allowed to cluster the sequential data up to 2 different maximum number of clusters; 1) from three up to five clusters (first row in Table 1), and 2) from three up to nine clusters (second row in Table 1). In order to overcome the problem of initialization of the K-means, the algorithm was run using 10 different initializations. For each clustering configuration, the DB-index (Bezdek et al, 1998) was used to measure the goodness of clustering. According to the DB-index measure, the number of states (clusters) in each model was determined according to the clustering configuration corresponding to the lowest value of the DB-index. Table 1 shows the number of states for each model for each clustering configuration.

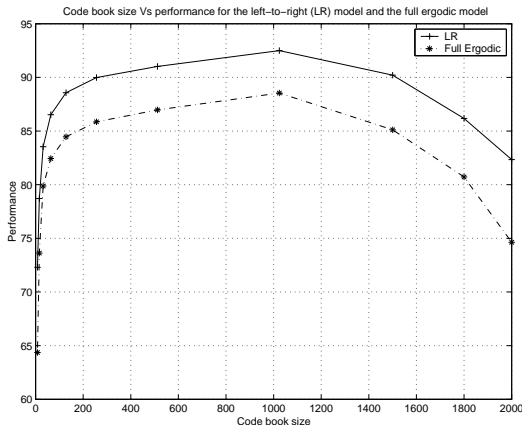


Fig. 4. Performance comparison between full ergodic and left-to-right models with different codebook sizes.

4.4 Studying the effect of model topology

To study the effect of the model topology on the performance, two HMM-based classifiers were considered. Both classifiers had the same number of models and the same number of states in each model but the model topology was different in both classifiers. The first classifier had full ergodic (fully connected) models while the second had left-to-right topology as described earlier. The experiment was conducted using ten code books (previous 7 plus 3 more with size 1500, 1800, 2000), five different initializations for the A matrix; (0.5 & 0.5, 0.6 & 0.4, 0.7 & 0.3, 0.8 & 0.2 and 0.9 & 0.1) for the ij and ii transitions respectively of the left-to-right model, and 5 different random initializations for the ergodic model. Figure 4 illustrates the results obtained from this experiment on the ten code books and the first initialization of each model.

As expected, the results show that the simpler model; which is the left-to-right in that case, always outperforms the full ergodic model. The full ergodic model represents a fully connected graph and hence has the largest number of parameters. According to the Bayesian approach, the model has the highest likelihood of the data which led the model to overfit the training set and hence the poor performance on the test set. As for the diffusion of credits factor, the A matrix for the full ergodic model does not have deterministic (0 or 1 probabilities) transitions which made it difficult for the model to learn long range dependencies.

The degradation of performance in Figure 4 is due to an accumulated effect from the vector quantization and the training of HMMs and it may be explained as follows. The vector quantization process was performed on the training set of the database and increasing the code book size led the algorithm to form smaller and finer (might be noise) clusters from the training set. Hence, the result is a well fitted code book for the training set and very

sensitive to slight variations, i.e over fitting. Next, the discrete HMMs used this sensitive but large code book for training, which implies that the HMMs were trained on very special sequences of symbols that may not occur in the test set. Consequently, the HMMs had over-fit the training set and will have a poor generalization on the test set. Hence, the fall in the two curves is due to the accumulated over-fit effect that started from the vector quantization and propagated to the HMM training.

5 Conclusion

We studied the effect of number of states and the topology on the performance of HMM-based classifiers. The Bayesian approach for model selection with the Ockham's Razor showed that simpler models will have better generalization than full ergodic (fully connected) models. On the other hand, to avoid any diffusion of credits while learning HMMs, transition probabilities should be deterministic (0 or 1 probabilities). Both of these results supported with the empirical experiments show that the topology has a stronger influence than the number of states in improving the modeling capability of HMMs and hence increasing the performance of HMM-based classifiers. It can be seen from Figures 2 and 4 that increasing the number of states from 3 to 6 increased the performance by almost 2% and changing the topology increased the performance by (4-5%). The result encourages us to design algorithms for HMMs, different than model selection techniques, that can learn the topology from the training data, i.e. set 0 or 1 transitions in the A matrix, especially in the absence of the a priori knowledge.

6 Acknowledgments

We would like to thank Incheol Kim from CENPARMI for useful discussions. Also, we would like to thank the Ministry of Education of Québec and NSERC of Canada for the financial support.

References

- L. E. Baum and T. Petrie, Statistical Inference for Probabilistic Functions of Finite State Markov Chains, *Ann. Math. Stat.* **37** (1966) 1554–1563.
- Y. Bengio, Markovian Models for Sequential Data, *Neural Computing Surveys* **41:1** (1999) 129–162.

- Y. Bengio and P. Frasconi, Diffusion of Credits in Markovain Model, *Neural Information Proccesing Systems* **7** (1995) 1251–1254.
- E. Senta, *Nonnegative Matrices and Markov Chains* (Springer, New York, 1986).
- L. R. Rabiner, A Tutorial on Hidden Markov Models and Selected Application in Speech Recognition, *Proc. IEEE* **77:2** (1989) 257–286.
- J. K. Baker, The Dragon system - An Overview, *IEEE Trans. Accoustics, Speech and Signal Proc* **23:11** (1975) 23–29.
- M. Bicego, A. Dovier and V. Murino, Designing the Minimal Structure Hidden Markov Model by Bisimulation, in: M. Figueiredo, J. Zerubia and A. K. Jain, eds., *Energy Minimization Methods in Computer Vision and Pattern Recognition* (Springer, 2001) 75–90.
- M. Bicego, V. Murino and M. Figueiredo, A Sequential Pruning Strategy for the Selection of the Number of States in Hidden Markov Models, *Pattern Recognition Letters* **24** (2003) 1395–1407.
- A. Stolcke and S. Omuhundro, Hidden Markov Model Induction by Bayesian Model Merging, in: S. Hanson, J. Cowan and C. Giles, eds., *Advances in Neural Information Processing 5* (Morgan Kaufmann, 1992) 11-18.
- T. Brants, Estimating Markov Model Structures, *Proc. of ICLSP Philadelphia, PA* (1996).
- A. El-Yacoubi, M. Gilloux, R. Sabourin and C. Y. Suen, An HMM Based Approach for Off-line Unconstrained Handwritten Word Modeling and Recognition, *IEEE Trans. PAMI* **21:8** (1999) 752–760.
- J. Lee, J. Kim and J-H Kim, Data-driven Design of HMM Topology for On-Line Handwritten Recognition, *Pattern Recognition* **15:1** (2001) 107–121.
- A. S. De Britto, R. Sabourin, F. Bortolozzi, and C. Y. Suen, An Enhanced HMM Topology in an LBA Framework for the Recognition of Handwritten Numeral Strings, *Int. Conf. on Advances in Pattern Recognition* (2001) 105–114.
- J. Lien, Automatic Recognition of Facial Expressions using Hidden Markov Models and Estimation of Expression Entensity, *CMU, School of Computer Science, Pittsburg* (PhD Thesis. 1998).
- R. Lyngso, C. Pedersen and H. Nielsen, Metrics and Similarity Measures for Hidden Markov Models, *Proc. Int. Conf. on Intelligent Systems for Molecular Biology* (1999) 178–186.
- C. Bahlman, H. Burkhardt and A. Ludwigs, Measuring HMM Similarity with the Bayes Probability of Error and its Application, *Proc. of 6th ICDAR, Seattle, Washington, USA* (2001) 406-411.
- V. Balasubramanian, Equivalence and Reduction of Hidden Markov Models, *MIT Aritifical Intelligence Laboratory Tech. Report* **1370** (1993).
- A. Biem, A Model Selection Criterion for Classification: Application to HMM Topology Optimization, *Proc. 17th ICDAR, Edinburgh, U.K* (2003) 104–108.
- E. Augusting, O. Baret, S. Knerr and D. Price, Hidden Markov Model Based Word Recognition and its Application to Legal Amount Recognition on

- Frensh Bank Cheques, *Computer Vision and Image Understanding* **70:3** (1998) 404–419.
- S. Cornell, A Comparison of Hidden Markov Model Features for the Recognition of Cursive Handwriting, *Dept. of Computer Science, Michigan State University* (Master Thesis, 1996)
- J. Cai and Z-Q. Liu, Hidden Markov Models with Spectral Features for 2D Shape Recognition, *IEEE Trans. PAMI* **23:12** (2001) 703–713.
- P. Bladi and S. Brunak, *Bioinformatics, The Machine Learning Approach* (MIT Press, 1998).
- K. Karplus, K. Sjolander, C. Barrett, M. Cline, D. Haussler, R. Hughey, L. Holm and C. Sander, Predicting Protein Structure Using Hidden Markov Models, *Proteins: Structure, Function and Genetics* **1:1** (1997) 134–139.
- I. Kim and S. Chien, Analysis of 3D Hand Trajectory Gestures Using Stroke-Based Composite Hidden Markov Models, *Applied Intelligence* **15** (2001) 131–143.
- R. Gray, Vector Quantization, *IEEE Trans. ASSP* (1984) 4–29.
- S. Theodoridis and K. Koutroubas, *Pattern Recognition* (Academic Press, 1999).
- R. O. Duda, P. E. Hart and D. G. Stork, *Pattern Classification, Second Edition* (Wiley-Interscience, 2001).
- J. C. Bezdek and N. R. Pal, Some New Indexes of Cluster Validity, *IEEE Trans. on Sys. Man and Cybernetics Part B* **28:3** (1998) 301–315
- D. Heckerman, A Tutorial on Learning with Graphical Models, MSR-TR-9506, (Microsoft Research, 1996).
- K. P. Murphy, A Introduction to Graphical Models, <http://www.ai.mit.edu/~murphyk/papers.html> (2001).
- Y. LeCun, The MNIST Database of Handwritten Digits, <http://yann.lecun.com/exdb/mnist>.
- C-L. Liu, K. Nakashima, H. Sako and H. Fujisawa, Handwritten Digit Recognition: Benchmarking of State-of-the-art Techniques, *Pattern Recognition* **36** (2003) 2271–2285.
- P. Simard and D. Steinkraus and J. C. Platt, Best Practices for Convolutional Neural Networks Applied to Visual Document Analysis, *Proc. of 17th IC-DAR, Edinburgh, U.K* (2003) 962–965.
- J. Dong, Speed and accuracy: Large-scale machine learning algorithms and their applications, *CENPARMI, Department of Computer Science, Concordia University, Montreal, Canada* (Ph.D Thesis, 2003).