



Available online at www.sciencedirect.com



Pattern Recognition
Letters

Pattern Recognition Letters xxx (2004) xxx–xxx

www.elsevier.com/locate/patrec

On the structure of hidden Markov models

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

^a Department of Computer Science, CENPARMI, Concordia University, GM-606, 1455 de Maisonneuve W., Montréal, QC, Canada H3G 1M8

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

Received 22 July 2003; received in revised form 8 February 2004

8 Abstract

9 This paper investigates the effect of HMM structure on the performance of HMM-based classifiers. The investi-
10 gation is based on the framework of graphical models, the diffusion of credits of HMMs and empirical experiments.
11 Although some researchers have focused on determining the number of states, this study shows that the topology has a
12 stronger influence on increasing the performance of HMM-based classifiers than the number of states.
13 © 2004 Published by Elsevier B.V.

14 *Keywords:* HMM structure; Graphical models; Credits diffusion; Ocham's razor; K-Means clustering

15 1. Introduction

16 Hidden Markov models (HMMs) (Baum and
17 Petrie, 1966) are a class of stochastic processes that
18 is capable of modeling time-series data. They be-
19 long to a larger class of models known as gener-
20 ative models. Though generative models are tools
21 for data modeling, in the literature, HMMs were
22 used in many classification problems such as
23 speech recognition (Rabiner, 1989; Baker, 1975),
24 handwritten word recognition (El-Yacoubi et al.,
25 1999), object recognition (Cai and Liu, 2001),
26 gesture recognition (Kim and Chien, 2001), bio-

informatics (Bladi and Brunak, 1998) and model- 27
ing of biological sequences (Karplus et al., 1997). 28

29 Designing an HMM for data modeling to be a
30 part of an HMM-based classifier means deter-
31 mining the structure (the number of states, and the
32 topology) of the model. The topology in this
33 context is meant to be the connections (transitions)
34 between the states. The structure affects the mod-
35 eling capability considerably and consequently the
36 performance of the classifier. An estimation of the
37 weight of each factor on the performance can
38 point out to the main factor affecting the perfor-
39 mance and consequently can lead to an improve-
40 ment in the selection of values of this factor. 40
41 Although some researchers focused on the prob-
42 lem of number of states and the topology, the goal
43 of this paper is to investigate the effect of the
44 number of states and the topology, each sepa-
45 rately, on the performance of HMM-based classi-

* Corresponding author. Tel.: +1-514-848-7953; fax: +1-514-848-4522.

E-mail address: k_aboumo@cenparmi.concordia.ca (K.T. Abou-Moustafa).

46 fiers. Our research shows that the topology can
47 improve the modeling capability greatly.

48 The investigation is based on (1) linking the
49 theoretical results from model selection for
50 graphical models with the diffusion of credits in
51 Markovian models and, (2) supporting the results
52 with empirical experiments for the recognition of
53 unconstrained handwritten digits. The experiments
54 compared the performance obtained from classi-
55 fiers with different structures.

56 It is worth calling the readers' attention to two
57 issues regarding this work. First, though HMMs
58 are usually trained with time-series data with a
59 long signal duration, such as the applications
60 mentioned above, the experiments used isolated
61 handwritten digits that have a short signal dura-
62 tion. The reason for that is to treat the HMM-
63 based classifier like other classifiers such as multi
64 layer perceptrons (MLPs) and support vector
65 machines (SVMs) without any constraints on the
66 data. Second, the goal of the paper is not to
67 introduce a new state-of-the-art recognition result
68 on the MNIST database using HMMs, but rather,
69 to compare the performance of HMM-based
70 classifiers under different structure conditions. It is
71 well known from the literature that classifiers such
72 as SVMs and MLPs can achieve very high recog-
73 nition results (Dong, 2003; Simard et al., 2003; Liu
74 et al., 2003) on this database.

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

82 The rest of the paper is organized as follows:
83 Section 2 reviews related work in the literature
84 and HMMs. Section 3 discusses the effect of the
85 structure on the modeling capability of HMMs.
86 Section 4 describes the experiments and results,
87 and finally Section 5 concludes the paper.

88 2. Related work

89 A work directly related to this investigation is
90 the problem of optimizing HMM structure in two

forms; (1) application dependent methods, and (2) 91
application independent methods. Application 92
dependent methods use a priori knowledge from 93
the application domain such as (El-Yacoubi et al., 94
1999), where they used information extracted from 95
a character segmentation process to build a special 96
HMM structure (number of states and the topol- 97
ogy), (De Britto et al., 2001) modified the left-to- 98
right model to enhance the performance of his 99
proposed framework for numeral strings and (Lee 100
et al., 2001) fixed the topology to be a left-to-right 101
one and determined the number of states by 102
reflecting the structure of a target pattern. The 103
major drawback of these methods is that they are 104
designed for specific applications and cannot be 105
generalized to others. 106

107 On the other hand, application independent 107
methods, although are more promising, yet they 108
are not popularized. These methods include the 109
work of (Stolcke and Omuhundro, 1992) and 110
(Brants, 1996) where they proposed an incremental 111
learning for the structure based on state merging 112
and splitting, i.e., the structure is changed as new 113
evidence is added to the model; (Lien, 1998) pro- 114
posed a general method to determine the number 115
of states and the connections between states in 116
discrete left-to-right HMMs. Recently, Bicego et 117
al. focused only on determining the number of 118
states using probabilistic bisimulation (Bicego et 119
al., 2001) and sequential pruning using Bayesian 120
information criterion (BIC) (Bicego et al., 2003). 121
Model selection approaches were also investigated 122
for this purpose and recently (Biem, 2003) pro- 123
posed a discriminative information criterion (DIC) 124
framework and used it to optimize the HMM 125
structure. 126

127 Different approaches for structure optimization 127
can also be found. (Lyngso et al., 1999) focused on 128
comparing HMMs in terms of state emission 129
probabilities, (Bahlman et al., 2001) used Bayesian 130
estimates of HMM states as a criterion for select- 131
ing HMMs and (Balasubramanian, 1993) selected 132
HMMs based on equal probabilities of observa- 133
tion sequences only. By examining the above lit- 134
erature, and except for the work Stolcke and 135
Brants, most of these methods use the left-to-right 136
topology and the optimization targets only the 137
number of states and the number of mixtures in 138

139 cases of continuous HMMs. However, according
140 to this investigation, the number of states may not
141 affect the performance after a certain limit, but it
142 can reduce the computational time for training
143 and testing, while the topology can considerably
144 affect the performance of HMMs.

145 3. HMM structure

146 In many applications that use HMMs, the
147 number of states is manually predetermined prior
148 to training. The connections between states,
149 (topology) is determined by setting non-zero
150 probabilities in the A matrix prior training. During
151 training, the EM (Baum–Welch) algorithm im-
152 proves the estimates of these probabilities from the
153 data. Note that the EM algorithm cannot set 0 or 1
154 (can approach 0 or 1) probabilities in the A matrix,
155 therefore it cannot be seen as an algorithm that
156 learns the topology. In the following, we investi-
157 gate the effect of the topology on the performance
158 of HMM-based classifiers through two different
159 perspectives: (1) using the graphical models
160 framework and, (2) using the diffusion of credits
161 while learning Markovian models.

162 3.1. Bayesian formulation for model selection

163 Determining the number of states and the
164 topology of HMMs can be viewed as a model
165 selection problem. The problem can be formulated
166 as follows. Given the training set of examples Ψ
167 and a criterion function \mathcal{T} for the quality of the
168 model on the data set Ψ , choose a model from a
169 certain set of models, in such a way to maximize
170 the expected value of this criterion function on
171 new data (assumed to be sampled from the same
172 unknown distribution from which the training
173 data was sampled) (Bengio, 1999).

174 HMMs can be viewed as a special case of
175 graphical models (Heckerman, 1996; Murphy,
176 2001). Model selection is one of the main problems
177 in graphical models and much work has been
178 introduced regarding this problem. The Bayesian
179 approach, one of the main approaches for model
180 selection, is a fundamental approach for model
181 selection in graphical models. Following this ap-

proach means encoding the uncertainty about the
182 structure of the HMM by using a discrete variable
183 whose states correspond to the possible HMM
184 structure hypotheses S^h and assessing it the a pri-
185 ori density $P(S^h)$. Given the training example set Ψ
186 for the model λ and augmenting the model
187 parameters A, B, π in a single parameter vector θ ,
188 the problem would be computing the posterior
189 distribution for the HMM structures. This can be
190 formulated as follows using Bayes theorem: 191

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

where $P(\Psi)$ is a constant that does not depend on
193 the network structure. 194

The maximum likelihood structure would be
195 the complete graph (Murphy, 2001), i.e., the full
196 ergodic model, since this has the greatest number
197 of parameters, and hence can achieve the highest
198 likelihood. On the other hand this increases the
199 model's complexity and will let the model overfit
200 the training data resulting in a poor generalization.
201 In fact, the marginal likelihood in 1 plays an
202 important role to prevent this overfit. From the
203 definition of the marginal likelihood: 204

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

it automatically penalizes more complex structures
206 since they have more parameters and hence cannot
207 give as much probability mass to the region of
208 space where the data actually lies. In other words,
209 a complex model is less believable and hence less
210 likely to be selected. This phenomenon is known as
211 Ockham's razor (Murphy, 2001) which favors sim-
212 ple models over complex ones. It can be seen that
213 though the number of states may be fixed, the
214 topology can affect the modeling capability in a
215 serious way. 216

3.2. Diffusion of credits in Markovian models 217

The work in (Bengio and Frasconi, 1995)
218 investigated the problem of diffusion in homoge-
219 neous and non-homogeneous HMMs and its effect
220 on learning long term dependencies. Training
221 HMMs requires propagating forward and back-
222 ward probabilities and taking products of the 223

224 transition matrix. Therefore, two types of diffusion
 225 exist, diffusion of influence in the forward path and
 226 diffusion of credit in the backward phase of
 227 training. The paper (Bengio and Frasconi, 1995)
 228 studied under which conditions these products of
 229 matrices will converge to a lower rank, thus
 230 harming learning long term dependencies. The
 231 difficulty of learning was measured by using the
 232 Dobrushin's ergodicity coefficient (Senta, 1986)
 233 defined as follows:

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

235 where $A = \{a_{ij}\}$ is the transition probability ma-
 236 trix. It was shown that in all cases, while training
 237 HMMs, the ergodicity coefficient will converge to
 238 0 indicating a greater difficulty in learning, but the
 239 rate of convergence depends on the topology. Fig.
 240 1 (Bengio and Frasconi, 1995) shows the conver-
 241 gence of four HMMs with the same number of
 242 states but with different topologies. It can be seen
 243 that the full ergodic model has the fastest conver-
 244 gence rate and that simpler models are slower.
 245 The final conclusion is that in order to avoid any
 246 kind of diffusion, most transition probabilities
 247 should be deterministic (0 or 1 probability). The
 248 result coincides with the Ocham's razor result
 249 obtained from the previous section and both prefer
 250 simple topologies over complicated ones.

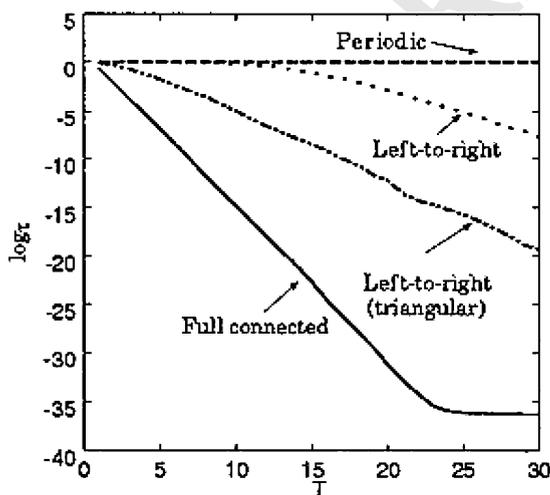


Fig. 1. Convergence of Dobrushin's coefficient for four different topologies.

4. Experiments

251

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

259

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 20×20 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

278

279

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 using the seven code books and with several initializations

293

294 for the A matrix as will be shown later. However,
295 due to the following reasons: (1) space limitation,
296 (2) avoid redundancy, and (3) similarity of results
297 and conclusions, we selected the clearest of these
298 experiments for illustration.

299 4.3. Studying the effect of number of states

300 In studying the effect of number of states, two
301 experiments were conducted. The first experiment
302 used HMMs with a left-to-right topology and all
303 models had an equal number of states. The
304 experiment studied the relation between the per-
305 formance and the increase in the number of states
306 in the classifier. The second experiment studied the
307 performance of classifiers with a varying number
308 of states in each model. It compared the perfor-
309 mance between models with an equal number of
310 states and models with a varying number of states.

311 4.3.1. Experiment 1

312 This experiment was conducted using the seven
313 code books, and for each experiment, the A matrix
314 was initialized using three different initializations;
315 (0.5 & 0.5, 0.7 & 0.3 and 0.9 & 0.1) for the ij and ii
316 transitions, respectively. Fig. 2 illustrates the re-
317 sults for *Experiment 1* using three code books and
318 the first initialization for the A matrix. It can be
319 seen that increasing the number of states can in-

crease the performance up to a certain limit, after
that a saturation is reached whenever more
unnecessary states are added. However, the satu-
ration 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 emission of symbols change. Let the true (un-
known) 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 unneces-
sary 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. Fig. 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).

Fig. 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 to large classification
problems. However, as followed in the literature
(El-Yacoubi et al., 1999; Augusting et al., 1998), a

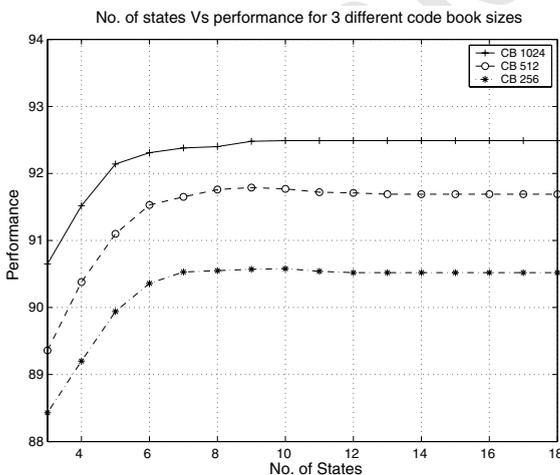


Fig. 2. The relation between performance and the number of states with different code book sizes.

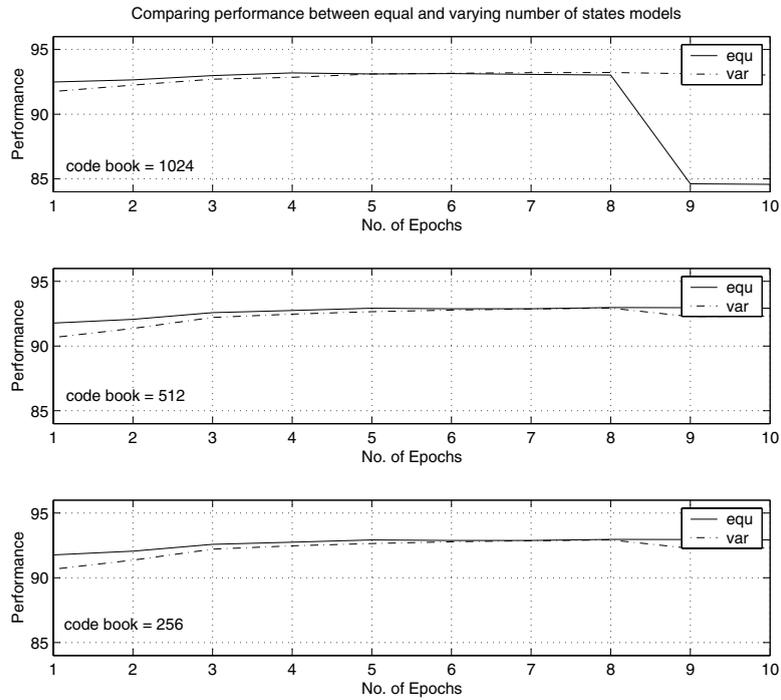


Fig. 3. Performance comparison between models with equal (EQU) and varying (VAR) number of states with different codebook sizes.

367 guaranteed performance with an easy design
 368 would be an HMM classifier with an equal number
 369 of states for all models. In Fig. 3, it is worth
 370 mentioning that the drop seen in the first graph is
 371 experienced in the other graphs for the EQU and
 372 VAR models but in late epochs not shown in the
 373 graphs. The reason for the drop is due to the
 374 overfit of models on the data and due to the dif-
 375 fusion of credits while learning.

376 4.3.3. Determining the number of states

377 As mentioned earlier, the number of states is
 378 usually fixed (manually predetermined). Excep-
 379 tions are models that use automatic clustering
 380 algorithms that determine the number of states
 381 and their outputs, but this still leaves out the
 382 topology (Brants, 1996; Theodoridis and Kou-

troumbas, 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
 and Pal, 1998) to measure the goodness of different
 clustering configurations and then select the best
 number of clusters according to this cluster valid-
 ity 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 two
 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
The number of states of each model

Model	0	1	2	3	4	5	6	7	8	9
No. of states (3–5)	5	5	5	5	4	4	4	5	3	4
No. of states (3–9)	6	5	8	8	9	6	8	8	3	9

400 Table 1). In order to overcome the problem of
401 initialization of the K-means, the algorithm was
402 run using 10 different initializations. For each
403 clustering configuration, the DB-index (Bezdek
404 and Pal, 1998) was used to measure the goodness
405 of clustering. According to the DB-index measure,
406 the number of states (clusters) in each model was
407 determined according to the clustering configura-
408 tion corresponding to the lowest value of the DB-
409 index. Table 1 shows the number of states for each
410 model for each clustering configuration.

411 4.4. Studying the effect of model topology

412 To study the effect of the model topology on the
413 performance, two HMM-based classifiers were
414 considered. Both classifiers had the same number
415 of models and the same number of states in each
416 model but the model topology was different in
417 both classifiers. The first classifier had full ergodic
418 (fully connected) models while the second had left-
419 to-right topology as described earlier. The experi-
420 ment was conducted using ten code books (previ-
421 ous 7 plus 3 more with size 1500, 1800, 2000), five
422 different initializations for the A matrix; (0.5 & 0.5,
423 0.6 & 0.4, 0.7 & 0.3, 0.8 & 0.2 and 0.9 & 0.1) for
424 the ij and ii transitions, respectively of the left-to-
425 right model, and five different random initializa-
426 tions for the ergodic model. Fig. 4 illustrates the

427 results obtained from this experiment on the ten
428 code books and the first initialization of each
429 model.

430 As expected, the results show that the simpler
431 model; which is the left-to-right in that case, al-
432 ways outperforms the full ergodic model. The full
433 ergodic model represents a fully connected graph
434 and hence has the largest number of parameters.
435 According to the Bayesian approach, the model
436 has the highest likelihood of the data which led the
437 model to overfit the training set and hence the
438 poor performance on the test set. As for the dif-
439 fusion of credits factor, the A matrix for the full
440 ergodic model does not have deterministic (0 or 1
441 probabilities) transitions which made it difficult for
442 the model to learn long range dependencies.

443 The degradation of performance in Fig. 4 is due
444 to an accumulated effect from the vector quanti-
445 zation and the training of HMMs and it may be
446 explained as follows. The vector quantization
447 process was performed on the training set of the
448 database and increasing the code book size led the
449 algorithm to form smaller and finer (might be
450 noise) clusters from the training set. Hence, the
451 result is a well fitted code book for the training set
452 and very sensitive to slight variations, i.e., over
453 fitting. Next, the discrete HMMs used this sensi-
454 tive but large code book for training, which im-
455 plies that the HMMs were trained on very special
456 sequences of symbols that may not occur in the
457 test set. Consequently, the HMMs had over-fit the
458 training set and will have a poor generalization on
459 the test set. Hence, the fall in the two curves is due
460 to the accumulated over-fit effect that started from
461 the vector quantization and propagated to the
462 HMM training.

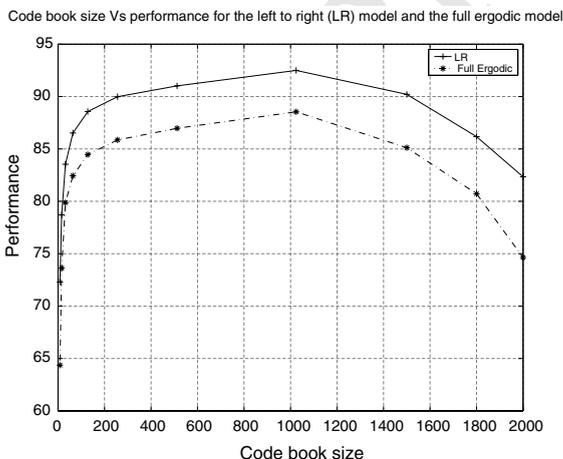


Fig. 4. Performance comparison between full ergodic and left-to-right models with different codebook sizes.

5. Conclusion

463 We studied the effect of number of states and
464 the topology on the performance of HMM-based
465 classifiers. The Bayesian approach for model
466 selection with the Ocham's razor showed that
467 simpler models will have better generalization than
468 full ergodic (fully connected) models. On the other
469 hand, to avoid any diffusion of credits while
470 learning HMMs, transition probabilities should be
471

472 deterministic (0 or 1 probabilities). Both of these
473 results supported with the empirical experiments
474 show that the topology has a stronger influence
475 than the number of states in improving the mod-
476 eling capability of HMMs and hence increasing the
477 performance of HMM-based classifiers. It can be
478 seen from Figs. 2 and 4 that increasing the number
479 of states from 3 to 6 increased the performance by
480 almost 2% and changing the topology increased
481 the performance by (4–5%). The result encourages
482 us to design algorithms for HMMs, different than
483 model selection techniques, that can learn the
484 topology from the training data, i.e., set 0 or 1
485 transitions in the A matrix, especially in the ab-
486 sence of the a priori knowledge.

487 Acknowledgements

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

493 References

494 Augusting, E., Baret, O., Knerr, S., Price, D., 1998. Hidden
495 Markov model based word recognition and its application
496 to legal amount recognition on French bank cheques.
497 *Comput. Vision Image Understand.* 70 (3), 404–419.
498 Bahlman, C., Burkhardt, H., Ludwigs, A., 2001. Measuring
499 HMM similarity with the Bayes probability of error and its
500 application. In: *Proc. 6th ICDAR*, Seattle, Washington,
501 USA. pp. 406–411.
502 Baker, J.K., 1975. The Dragon system—an overview. *IEEE*
503 *Trans. Acoust., Speech Signal Process* 23 (11), 23–29.
504 Balasubramanian, V., 1993. Equivalence and reduction of
505 hidden Markov models. *MIT Artificial Intell. Lab. Tech.*
506 *Report 1370*.
507 Baum, L.E., Petrie, T., 1966. Statistical inference for probabi-
508 listic functions of finite state Markov chains. *Ann. Math.*
509 *Statist.* 37, 1554–1563.
510 Bengio, Y., Frasconi, P., 1995. Diffusion of credits in Markov-
511 vain model. *Neural Informat. Process. System* 7, 1251–1254.
512 Bengio, Y., 1999. Markovian models for sequential data.
513 *Neural Comput. Surveys* 41 (1), 129–162.
514 Bezdek, J.C., Pal, N.R., 1998. Some new indexes of cluster
515 validity. *IEEE Trans. Systems Man Cybernet. Part B* 28 (3),
516 301–315.

Bicego, M., Dovier, A., Murino, V., 2001. Designing the
517 minimal structure hidden Markov model by bisimulation.
518 In: *Figueiredo, M., Zerubia, J., Jain, A.K. (Eds.), Energy*
519 *Minimization Methods in Computer Vision and Pattern*
520 *Recognition*. Springer, pp. 75–90.
521 Bicego, M., Murino, V., Figueiredo, M., 2003. A sequential
522 pruning strategy for the selection of the number of states in
523 hidden Markov models. *Pattern Recognition Lett.* 24, 1395–
524 1407.
525 Biem, A., 2003. A model selection criterion for classification:
526 Application to HMM topology optimization. In: *Proc. 17th*
527 *ICDAR*, Edinburgh, UK. pp. 104–108.
528 Bladi, P., Brunak, S., 1998. *Bioinformatics, the Machine*
529 *Learning Approach*. MIT Press.
530 Brants, T., 1996. Estimating Markov model structures. In:
531 *Proc. ICLSP Philadelphia, PA*.
532 Cai, J., Liu, Z.-Q., 2001. Hidden Markov models with spectral
533 features for 2D shape recognition. *IEEE Trans. PAMI* 23
534 (12), 703–713.
535 Cornell, S., 1996. A comparison of hidden Markov model
536 features for the recognition of cursive handwriting. Depart-
537 ment of Computer Science, Michigan State University,
538 Master Thesis, 1996.
539 De Britto, A.S., Sabourin, R., Bortolozzi, F., Suen, C.Y., 2001.
540 An enhanced HMM topology in an LBA framework for the
541 recognition of handwritten numeral strings. In: *Internat.*
542 *Conf. on Advances in Pattern Recognition*. pp. 105–114.
543 Dong, J., 2003. Speed and accuracy: Large-scale machine
544 learning algorithms and their applications, CENPARMI,
545 Department of Computer Science, Concordia University,
546 Montreal, Canada, Ph.D. Thesis, 2003.
547 Duda, R.O., Hart, P.E., Stork, D.G., 2001. *Pattern Classifica-*
548 *tion*, second ed. Wiley-Interscience.
549 El-Yacoubi, A., Gilloux, M., Sabourin, R., Suen, C.Y., 1999.
550 An HMM based approach for off-line unconstrained
551 handwritten word modeling and recognition. *IEEE Trans.*
552 *PAMI* 21 (8), 752–760.
553 Gray, R., 1984. Vector quantization. *IEEE Trans. ASSP*, 4–29.
554 Heckerman, D., 1996. A tutorial on learning with graphical
555 models, MSR-TR-9506, Microsoft Research.
556 Karplus, K., Sjolander, K., Barrett, C., Cline, M., Haussler, D.,
557 Hughey, R., Holm, L., Sander, C., 1997. Predicting protein
558 structure using hidden Markov models. *Proteins: Struct.,*
559 *Funct. Genet.* 1 (1), 134–139.
560 Kim, I., Chien, S., 2001. Analysis of 3D hand trajectory
561 gestures using stroke-based composite hidden Markov
562 models. *Appl. Intell.* 15, 131–143.
563 LeCun, Y., 1998. The MNIST Database of Handwritten Digits,
564 yann.lecun.com/exdb/mnist.
565 Lee, J., Kim, J., Kim, J.-H., 2001. Data-driven design of HMM
566 topology for on-line handwritten recognition. *Pattern Rec-*
567 *ognition* 15 (1), 107–121.
568 Lien, J., 1998. Automatic recognition of facial expressions using
569 hidden Markov models and estimation of expression inten-
570 sity. CMU, School of Computer Science, Pittsburg, Ph.D.
571 Thesis, 1998.
572

- 573 Liu, C.-L., Nakashima, K., Sako, H., Fujisawa, H., 2003. 585
574 Handwritten digit recognition: Benchmarking of state-of- 586
575 the-art techniques. *Pattern Recognition* 36, 2271–2285. 587
576 Lyngso, R., Pedersen, C., Nielsen, H., 1999. Metrics and 588
577 similarity measures for hidden Markov models. In: Proc. 589
578 Internat. Conf. on Intelligent Systems for Molecular Biol- 590
579 ogy. pp. 178–186. 591
580 Murphy, K.P., 2001. A Introduction to Graphical Models, 592
581 www.ai.mit.edu/murphyk/papers.html. 593
582 Rabiner, L.R., 1989. A tutorial on hidden Markov models and 594
583 selected application in speech recognition. *Proc. IEEE* 77 595
584 (2), 257–286. 596
- Senta, E., 1986. *Nonnegative matrices and Markov chains*. 585
Springer, New York. 586
Simard, P., Steinkraus, D., Platt, J.C., 2003. Best practices for 587
convolutional neural networks applied to visual document 588
analysis. In: Proc. 17th ICDAR, Edinburgh, UK. pp. 962– 589
965. 590
Stolcke, A., Omuhundro, S., 1992. Hidden Markov model 591
induction by Bayesian model merging. In: Hanson, S., 592
Cowan, J., Giles, C. (Eds.), *Advances in Neural Informa- 593
tion Processing 5*. Morgan Kaufmann, pp. 11–18. 594
Theodoridis, S., Koutroumbas, K., 1999. *Pattern Recognition*. 595
Academic Press. 596