

Document downloaded from:

<http://hdl.handle.net/10251/50978>

This paper must be cited as:

Giménez Pastor, A.; Andrés Ferrer, J.; Juan, A. (2014). Discriminative Bernoulli HMMs for isolated handwritten word recognition. *Pattern Recognition Letters*. 35:157-168.
doi:10.1016/j.patrec.2013.05.016.



The final publication is available at

<http://dx.doi.org/10.1016/10.1016/j.patrec.2013.05.016>

Copyright Elsevier

Discriminative Bernoulli HMMs for isolated handwritten word recognition

Adrià Giménez, Jesús Andrés-Ferrer and Alfons Juan

DSIC, Universitat Politècnica de València, Camí de Vera s/n, 46022 València, Spain

Abstract

Bernoulli HMMs (BHMMs) have been successfully applied to handwritten text recognition (HTR) tasks such as continuous and isolated handwritten words. BHMMs belong to the generative model family and, hence, are usually trained by (joint) maximum likelihood estimation (MLE) by means of the Baum-Welch algorithm. Despite the good properties of the MLE criterion, there are better training criteria such as maximum mutual information (MMI). The MMI is the most widespread criterion to train discriminative models such as log-linear (or maximum entropy) models. Inspired by a BHMM classifier, in this work, a log-linear HMM (LLHMM) for binary data is proposed. The proposed model is proved to be equivalent to the BHMM classifier, and, in this way, a discriminative training framework for BHMM classifiers is defined. The behavior of the proposed discriminative training framework is deeply studied in a well known task of isolated word recognition, the RIMES database.

Keywords: HTR, Bernoulli HMM, Log-Linear HMM, MMI, RIMES

Email address: {agimenez,jandres,ajuan}@dsic.upv.es (Adrià Giménez, Jesús Andrés-Ferrer and Alfons Juan)

Preprint submitted to Pattern Recognition Letters

March 24, 2014

1. Introduction

In the past few years Bernoulli HMMs (BHMMs) have been proved to be competitive for handwritten text recognition (HTR). Specifically, competitive performance has been reported by BHMMs on handwritten English text (Giménez and Juan, 2009), and Arabic HTR (Giménez et al., 2010; Märgner and El Abed, 2010)¹.

Handwritten word classifiers based on HMMs, and in particular in BHMMs, are *generative models*. Generative models are classifiers based on the optimal Bayes classifier (Duda and Hart, 1973) which classify choosing the class c^* that maximizes the posterior class probability for a given input \mathbf{x} as follows

$$c^*(\mathbf{x}) = \arg \max_c p(c | \mathbf{x}) = \arg \max_c p(c, \mathbf{x}) \quad (1)$$

where instead of directly approximating the posterior class probability $p(c | \mathbf{x})$, the joint probability is modelled by a distribution $p_\theta(c, \mathbf{x})$ parameterized with θ . Among many other advantages, the generative models have two outstanding properties. On the one hand, the parameters of the generative models are easily understandable for researchers. For instance, in the BHMMs all parameters can be interpreted as percentages, and in particular, the emission parameters that model the emission probabilities are easily interpreted as grey level images. On the other hand, generative models are mostly trained with maximum likelihood estimation (MLE) criterion for which there are well-known algorithms for training latent variable mod-

¹the BHMM achieved the first place prize in the Arabic HTR competition organized during the ICFHR 2010 conference.

22 els (Dempster et al., 1977) in general, and HMMs in particular (Rabiner and
23 Juang, 1993).

24 Despite the good properties of the MLE criterion, it acknowledges an
25 important drawback in classification problems. The MLE is aimed at ex-
26 plaining the probability distribution that underlies in the training sample by
27 maximizing the likelihood of the joint probability function $p_{\theta}(c, \mathbf{x})$. However,
28 we are simply interested in classifying inputs, and there is no guarantee that
29 the MLE parameters are the most suitable for classifying, even though they
30 have been proved to be competitive.

31 An alternative to generative models are the discriminative models. Dis-
32 criminative models and criteria are aimed at classifying the data without
33 explaining the data distribution itself. These models are also based on the
34 Bayes decision rule in (1) but instead of the joint probability, they directly
35 approximate the posterior class probability by a model $p_{\lambda}(c|\mathbf{x})$ parameter-
36 ized by λ . However, discriminative parameters are difficult to understand
37 provided that they do not explain the input. Discriminative parameters
38 are usually estimated by the *maximum mutual information (MMI)* criterion,
39 which directly maximizes the likelihood of the posterior probability function
40 $p_{\lambda}(c|\mathbf{x})$. In contrast to MLE, the parameters estimated with MMI maximize
41 the most the differences between classes in order to better classify samples.
42 Unfortunately, there is no closed form solution for the MMI criterion, and few
43 unsatisfactory algorithms are available for finding the optimal parameters.
44 This problem is specially remarkable for discriminative models with hidden
45 variables as HMMs.

46 In Giménez et al. (2011) a MMI training scheme for Bernoulli mixture

47 classifiers was proposed and tested in a task of isolated handwritten digit
48 recognition. The proposed approach was based on the idea of finding a simi-
49 lar discriminative classifier to the Bernoulli mixture classifier, and then prove
50 the equivalence between both classifiers. The results analyzed in Giménez
51 et al. (2011) report that discriminatively trained Bernoulli mixture classi-
52 fier outperforms the generative Bernoulli mixture classifiers. In this paper
53 the work in Giménez et al. (2011) is extended to more complex models, the
54 BHMMs, which are assessed in a complex isolated word recognition task.
55 Specifically, we compared both generative and discriminative approaches in
56 the RIMES database (Grosicki et al., 2009) in which the discriminative mod-
57 els obtained very competitive results surpassing the performance obtained
58 by generative classifiers.

59 More precisely, the contributions of this work are the following:

- 60 1. We propose a particular case of log-linear HMM (LLHMM) classifier,
61 which can also be interpreted as a semi-Markov conditional Markov
62 chain (semi-CRF), for binary data inspired by the BHMM classifier.
- 63 2. We prove the equivalence between BHMMs and the proposed discrim-
64 inative model for binary data.
- 65 3. We provide a discriminative training scheme for BHMM classifiers by
66 means of their equivalence with LLHMMs, and analyze several discrim-
67 inative training criteria such as MMI.
- 68 4. We provide the capability to understand discriminative parameters
69 from a generative point of view by means of their equivalence with
70 BHMMs.

71 The remainder of the paper is organized as follows. A review of BHMMs

72 is given in Sec. 2. The proposed LLHMM or semi-CRF classifier for binary
 73 data is described in Sec. 3. The Sec. 4 proves equivalence between both clas-
 74 sifiers, and in Sec. 5 the parameter estimation for the LLHMM is analyzed.
 75 The proposed training scheme is deeply analyzed on the RIMES database
 76 in Sec. 6. We conclude the paper by summarizing and discussing the most
 77 important results and future research directions.

78 2. Bernoulli HMM

79 Let $O = (\mathbf{o}_1, \dots, \mathbf{o}_T)$ be a sequence of feature vectors. An HMM is a
 80 probability (density) function of the form

$$p_\theta(O) := \sum_{\mathbf{q}} p_\theta(O, \mathbf{q}) = \sum_{\mathbf{q}} \prod_{t=0}^T a(q_t, q_{t+1}) \prod_{t=1}^T b_{q_t}(o_t), \quad (2)$$

81 where we have uncovered the latent variables $\mathbf{q} = (q_0, q_1, \dots, q_{T+1})$ which
 82 represent all the possible state sequences (or paths), such that $q_1, \dots, q_T \in$
 83 $\{1, \dots, M\}$ are the regular states chosen out of a total of M states, and the
 84 first ($q_0 = I$) and last ($q_{T+1} = F$) states are special, the so-called *initial or*
 85 *start state* and the *final or stop* state, respectively. Furthermore, for any
 86 regular states i and j , $a(i, j)$ denotes the *transition* probability from i to j ,
 87 while b_j is the probability distribution for an *observation* at state j .

88 If we further assume that O is a sequence of binary featured vectors, then
 89 a *Bernoulli (mixture) HMM (BHMM)* is an HMM in which the probability of
 90 observing \mathbf{o}_t at position t and the state j ($q_t = j$) follows a Bernoulli mixture
 91 distribution

$$b_j(\mathbf{o}_t) := \sum_{k=1}^K \tau_j(k) \prod_{d=1}^D p_{jkd}^{o_{td}} (1 - p_{jkd})^{1-o_{td}}, \quad (3)$$

92 where $\tau_j(k)$ and \mathbf{p}_{jk} are, respectively, the prior and prototype of the k -th
 93 mixture component in state j . Fig. 1, depicts some prototypes for several
 94 states and components.

95 In isolated handwriting word recognition, BHMMs are used in Bayes'
 96 classifiers to approximate the input probability of a binary image, which is
 97 represented by an observation sequence of binary feature vectors for a given
 98 transcription. More precisely, the most probable transcription $S^* \in W$ for a
 99 given observation sequence O is obtained according to

$$S^* = \arg \max_S p(S, O) = \arg \max_S p(S) p(O | S), \quad (4)$$

100 where for each possible transcription S , the emission probability, $p(O | S)$, is
 101 approximated as a BHMM, and $p(S)$ is modelled regarding each probability
 102 as a parameter itself, π_S .

103 The number of possible transcriptions in handwriting recognition is typi-
 104 cally large, and consequently, the resulting parameters for a BHMMs are very
 105 sparse. In order to amend sparseness problems, BHMMs at word level are
 106 built from shared, embedded BHMMs at character level. More precisely, let
 107 C be the number of different characters (symbols) from which global BHMMs
 108 are built, we assume that each character c is modeled with a different BHMM
 109 with parametric vector θ_c , which is shared among words. Unfortunately, the
 110 input featured vectors are not aligned with individual characters, and the
 111 character boundaries are, therefore, unknown. For a given feature sequence
 112 $O = (\mathbf{o}_1, \dots, \mathbf{o}_T)$ representing a sequence of symbols $S = (s_1, \dots, s_L)$, with
 113 $L \leq T$; the latent segmentation $\mathbf{i} = (i_1, i_2, \dots, i_{L+1})$ defined as follows

$$1 = i_1 < \dots < i_L < i_{L+1} = T + 1; \quad (5)$$

114 induces the segmentation of O into L segments which monotonically corre-
 115 spond to each symbol. Specifically, the feature segment corresponding to
 116 the l -th character, s_l , is denoted by $O(i_l, i_{l+1}) = \mathbf{o}_{i_l}, \dots, \mathbf{o}_{i_{l+1}-1}$. Finally, the
 117 probability of O is determined by

$$p_{\theta}(O | S) = \sum_{\mathbf{i}} p_{\theta}(O, \mathbf{i} | S) = \sum_{\mathbf{i}} \prod_{l=1}^L p(O(i_l, i_{l+1}) | s_l), \quad (6)$$

118 where the sum is carried out over all possible segmentations of O into L
 119 segments, and $p(O(i_l, i_{l+1}) | s_l)$ is the probability of the l -th segment given
 120 by a BHMM in (2) with the parameters, θ_{s_l} , associated to the character
 121 s_l . Note that θ comprises all the embedded character parameters, i.e. $\theta =$
 122 $\{\theta_1, \dots, \theta_C\}$. These parameters are commonly estimated by MLE (Giménez
 123 et al., 2010).

124 Many of the parameters of the discussed model are constrained to sum
 125 1, since they directly approximate probabilities. These parameters are the
 126 mixture coefficients, the state transitions and the word prior probabilities,
 127 which must verify

$$\sum_w \pi_w = 1, \quad \forall_{c,q} : \sum_{q'} a_c(q, q') = 1, \quad \forall_{c,q} : \sum_k \tau_{cq}(k) = 1. \quad (7)$$

128 In Fig. 1, an embedded BHMM for the number 313 is shown. Note that
 129 the word model is the result of concatenating character models for the digit
 130 3, blank space, digit 1, and the blank space and digit 3 again, in that order.

131 3. Log-linear HMM Classifier for Binary Data

132 In this section, we propose a discriminative classifier inspired by the
 133 BHMM classifier for isolated handwritten words (4). The discriminative

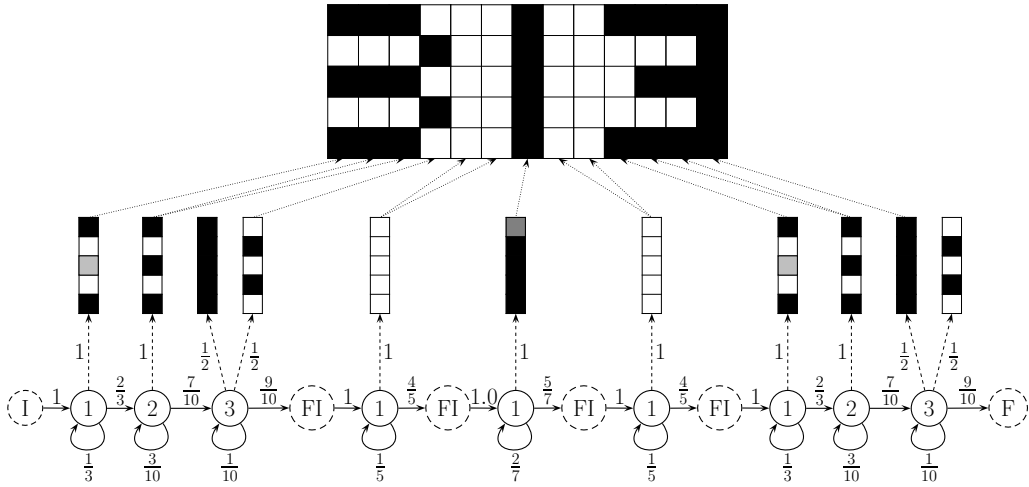


Figure 1: The binary featured vector representing the number 313 and the most probable path for generating it accordingly to a BHMM. Bernoulli prototype probabilities are represented using the following color scheme: black=1, white=0, gray=0.5 and light gray=0.1.

134 classifier proposed is based on a log-linear model, which is inferred from the
 135 parameters of a BHMM classifier. In what follows, we define the log-linear
 136 model and how a log-linear HMM (LLHMM) discriminative classifier can be
 137 built using it.

138 3.1. BHMM Inspired Log-linear Model

139 The BHMM classifier can be expressed by plugging (2), (3), and (6) as
 140 follows

$$p_{\theta}(O, S) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} p_{\theta}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}), \quad (8)$$

141 where by $\mathbf{i}, \mathbf{q}, \mathbf{k}$ we denote the 3 latent variables of the model, namely: the
 142 segmentation, \mathbf{i} , of O into L segments as defined in (5); the state sequence,
 143 $\mathbf{q} = (q_0, q_1, \dots, q_{T+1})$; and the emission component at each state, \mathbf{k} . Accord-

144 ing to the given segmentation \mathbf{i} , the state sequence \mathbf{q} must be valid, which
 145 implies that if t belongs to the l -th segment, then the state q_t must be a
 146 possible state of the character-level BHMM for the corresponding symbol s_l .
 147 Similarly, $\mathbf{k} = (k_1, \dots, k_T)$ must be a valid integer sequence where k_t denotes
 148 the selected mixture component for state q_t , among all the components of
 149 the state.

150 The joint probability in the right-hand-side of previous equation, $p_\theta(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})$,
 151 is decomposed left-to-right as follows

$$p_\theta(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \pi_S p_\theta(O, \mathbf{i}, \mathbf{q}, \mathbf{k} | S) = \pi_S p_\theta(\mathbf{i}, \mathbf{q} | S) p_\theta(O, \mathbf{k} | \mathbf{i}, \mathbf{q}, S) \quad (9)$$

152 where $p_\theta(\mathbf{i}, \mathbf{q} | S)$ is the transition probability of the word-level BHMM and
 153 $p_\theta(O, \mathbf{k} | \mathbf{i}, \mathbf{q}, S)$ the emission probability. The transition probabilities are
 154 then decomposed into

$$p_\theta(\mathbf{i}, \mathbf{q} | S) := \prod_{l=1}^L a_{s_l}(I, q_{i_l}) \cdot a_{s_l}(q_{i_{l+1}-1}, F) \prod_{t=i_l}^{i_{l+1}-2} a_{s_l}(q_t, q_{t+1}) \quad (10)$$

155 where the first product accounts for the input, $a_{s_l}(I, q_{i_l})$, and output, $a_{s_l}(q_{i_{l+1}-1}, F)$,
 156 transitions of the embedded model for the character s_l ; and where the sec-
 157 ond product are the inner transitions within the embedded character model.
 158 In the remaining of the paper, we will not differentiate between inner and
 159 outer transitions since this is a well-known characteristic of HMM, and by
 160 extension to our BHMM model. Furthermore, this significantly simplifies the
 161 notation. For instance, previous equation is expressed as

$$p_\theta(\mathbf{i}, \mathbf{q} | S) := \prod_{l,t} a_{s_l}(q_t q_{t+1}) \quad (11)$$

162 where we have also omitted the boundaries of the products, which can always
 163 be traced back to (10).

164 Similarly, the emission probability is decomposed as follows

$$p_\theta(O, \mathbf{k} \mid \mathbf{i}, \mathbf{q}, S) := \prod_{l=1}^L \prod_{t=i_l}^{i_{l+1}-1} \tau_{s_l q_t}(k_t) \prod_{d=1}^D p_{s_l q_t k_t d}^{o_{td}} (1 - p_{s_l q_t k_t d})^{(1-o_{td})}. \quad (12)$$

165 where again by dropping the product boundaries is simplified to

$$p_\theta(O, \mathbf{k} \mid \mathbf{i}, \mathbf{q}, S) := \prod_{l,t} \tau_{s_l q_t}(k_t) \prod_d p_{s_l q_t k_t d}^{o_{td}} (1 - p_{s_l q_t k_t d})^{(1-o_{td})}. \quad (13)$$

166 with $\tau_{s_l q_t}(k_t)$ and $\mathbf{p}_{s_l q_t k_t}$ being the prior and prototype of the k -th mixture
167 component at state q_t of the character s_l .

168 Consequently, the model in (9) can be expressed as follows

$$p_\theta(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \exp(\log \pi_S + \log p_\theta(\mathbf{i}, \mathbf{q} \mid S) + \log p_\theta(O, \mathbf{k} \mid \mathbf{i}, \mathbf{q}, S)) \quad (14)$$

169 where the logarithms of the probabilities are given by

$$\log p_\theta(\mathbf{i}, \mathbf{q} \mid S) = \sum_{l,t} \log a_{s_l}(q_t, q_{t+1}) \quad (15)$$

170 and

$$\log p_\theta(O, \mathbf{k} \mid \mathbf{i}, \mathbf{q}, S) = \sum_{l,t} (\log \tau_{s_l q_t}(k_t) + \xi_{s_l q_t}(k_t) + \sum_{l,t,d} o_{td} \log \frac{p_{s_l q_t k_t d}}{(1 - p_{s_l q_t k_t d})}), \quad (16)$$

171 with $\xi_{c q}(k)$ defined as

$$\xi_{c q}(k) = \sum_d \log(1 - p_{c q k d}). \quad (17)$$

172 Note that the term $\xi_{c q}(k)$ is easily obtained when applying the logarithm to
173 (13) by rearranging terms similarly to Giménez et al. (2011).

At this point, we reparameterize the probabilities in terms of the new parameters, λ , as follow

$$\lambda_S = \log \pi_S, \quad (18)$$

$$\lambda_{cqq'} = \log a_c(q, q'), \quad (19)$$

$$\lambda_{cqqk} = \log \tau_{cqq}(k) + \xi_{cqq}(k), \quad (20)$$

$$\lambda_{cqqkd} = \log \frac{p_{cqqkd}}{1 - p_{cqqkd}}, \quad (21)$$

174 for each character, c ; states, q and q' ; mixture component, k ; and input
175 dimension, d .

176 Provided the previous parameterization, the original joint probability
177 in (9) is alternatively expressed as follows

$$p_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \exp(\lambda_S + \sum_{l,t} \lambda_{s_l q_t q_{t+1}} + \sum_{l,t} \lambda_{s_l q_t k_t} + \sum_{l,t,d} o_{td} \lambda_{s_l q_t k_t d}) \quad (22)$$

178 In order to simplify notation, we adopt here the standard and powerful nota-
179 tion of log-linear models. We define an index m that ranges through all the
180 subindexes of the previous equation, i.e., m ranges from $\{S\}$ over $\{c, q, q'\}$
181 and $\{c, q, k\}$ to $\{c, q, k, d\}$. We also introduce a function $g_m(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})$ that
182 equals to the number of times the parameter λ_m is used, except for the pa-
183 rameters $\{\lambda_{cqqkd}\}$. In this case, the function $g_m(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})$ with $m = cqqkd$,
184 counts the number of times the d -th bit is *set* and has been generated with
185 the k -th component in the state, q , of the character, c . The simplest case
186 of the function is that of the prior parameters λ_S for which $g_m = 1$ (with
187 $m = S$).

188 The proposed notation simplifies (22) into

$$p_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \exp\left(\sum_{m \in \mathcal{M}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})} \lambda_m g_m(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})\right) \quad (23)$$

189 where by $\mathcal{M}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})$ we denote the set of values through which the index
190 m ranges. It is important to notice that this set depends on all the variables,
191 namely $O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}$; and changes with them. However, it is simpler to define
192 \mathcal{M} as the union of all the possible indexes that our parameters require and
193 replace the functions g_m by the so-called *feature functions* $f_m(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})$,
194 which are equal to g_m if m is an index of a required parameter and 0 otherwise.
195 For instance, consider again the word prior example with the new domain
196 \mathcal{M} . In this case, the index m can take the value of any word, S' in the
197 vocabulary; and then the feature function is defined as

$$f_{S'}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \delta(S, S') \quad (24)$$

198 where $\delta(a, b)$ is the Kronecker delta function, which equals 1 if both ele-
199 ments are equal, and 0 otherwise. The feature functions for the remaining
200 parameters are detailed in Sec. 3.3.

201 Finally, equation (22) is expressed as

$$p_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \exp\left(\sum_{m \in \mathcal{M}} \lambda_m f_m(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})\right) = \exp(\lambda' \mathbf{f}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})), \quad (25)$$

202 where we can substitute the sum by its vectorial notation. The model in (25)
203 when plugged into (8) is a log-linear model with binary inputs.

204 3.2. Discriminative Classifier

205 Log-linear models are commonly employed to approximate posterior prob-
206 abilities. From (25), we can approximate the posterior class probability re-
207 quired by the optimal Bayes' classifier in (1) as follows

$$p_\lambda(S | O) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} p_\lambda(S, \mathbf{i}, \mathbf{q}, \mathbf{k} | O), \quad (26)$$

208 where $p_\lambda(S, \mathbf{i}, \mathbf{q}, \mathbf{k} | O)$ is approximated by (25) and the Bayes' theorem as

$$p_\lambda(S, \mathbf{i}, \mathbf{q}, \mathbf{k} | O) = \frac{p_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})}{p_\lambda(O)} = \frac{\exp(\sum_{m \in \mathcal{M}} \lambda_m f_m(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}))}{p_\lambda(O)}. \quad (27)$$

209 It is worth noting that the denominator is a probability because of the
 210 transformation that we have performed in equations (18)-(21). However, we
 211 wish to select any arbitrary parametric vector, λ , and in such a case, the
 212 denominator also becomes arbitrary, yielding the *LLHMM model*

$$p_\lambda(S, \mathbf{i}, \mathbf{q}, \mathbf{k} | O) = \frac{\exp(\sum_{m \in \mathcal{M}} \lambda_m f_m(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}))}{\mathcal{Z}_\lambda(O)}, \quad (28)$$

213 where $\mathcal{Z}_\lambda(O)$ is a normalization constant defined as

$$\mathcal{Z}_\lambda(O) = \sum_S \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} \exp(\lambda' \mathbf{f}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})), \quad (29)$$

214 that for the specific parameters in equations (18)-(21) corresponds to the
 215 marginal probability $p_\lambda(O)$. The log-linear model in (28) is a log-linear model
 216 with hidden variables for the segmentation and for the states which have a
 217 first order dependence. This model is variation of a semi-Markov conditional
 218 random field.

219 The previous LLHMM is used in the optimal Bayes' rule to obtain the
 220 *LLHMM classifier*

$$S^* = \arg \max_S p_\lambda(S | O), \quad (30)$$

221 3.3. Feature Functions

222 As discussed before, in order to use the standard notation in log-linear
 223 models, we need to define the feature functions for each kind of parameters.

224 For a given character c out of C different symbols, and for a given pair
 225 of state indexes (q, q') of that character, we define the *transition features*
 226 $f_{cqq'}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = f_{cqq'}(S, \mathbf{i}, \mathbf{q})$ as follows

$$f_{cqq'}(S, \mathbf{i}, \mathbf{q}) = \sum_{l=1}^L \delta(s_l, c) \begin{cases} \sum_{t=i_l}^{i_{l+1}-2} \delta(q_t, q) \delta(q_{t+1}, q') & 1 \leq q, q' \leq M_c \\ \delta(q_{i_l}, q') & q = I, 1 \leq q' \leq M_c \\ \delta(q_{i_{l+1}-1}, q) & 1 \leq q \leq M_c, q' = F \\ 0 & \text{otherwise} \end{cases} \quad (31)$$

227 where I and F represent the *initial* and *final* states respectively. Intuitively,
 228 this feature counts the number of times the specific transition from q to q'
 229 of the character c , is used in the input $S, \mathbf{i}, \mathbf{q}$. Note that it can be 0, if, for
 230 instance, the character c is not part of word S .

231 For the mixture components, we define the *component features* for each
 232 character c , state q and component k as follows

$$f_{cqk}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = f_{cqk}(S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \sum_{l=1}^L \delta(s_l, c) \sum_{t=i_l}^{i_{l+1}-1} \delta(q_t, q) \delta(k_t, k). \quad (32)$$

233 Intuitively, this feature counts the number of times an emission of O is gen-
 234 erated by the k -th component of the state q of the character c .

235 The final set of features are the *emission features*, which are given as
 236 follows for each character c , state q , component k and dimension d

$$f_{cqkd}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \sum_{l=1}^L \delta(s_l, c) \sum_{t=i_l}^{i_{l+1}-1} \delta(q_t, q) \delta(k_t, k) o_{td}. \quad (33)$$

237 **4. Equivalence Between BHMMs and LLHMMs**

238 In this section we prove that the BHMM classifier for isolated words is
 239 equivalent to the LLHMM proposed in Sec. 3. A generative classifier is said
 240 to be equivalent to a discriminative classifier if for a given set of generative
 241 parameters θ , a set of discriminative parameters λ can be found such that

$$\arg \max_{S \in W} p_{\theta}(O, S) = \arg \max_{S \in W} p_{\lambda}(S | O); \quad (34)$$

242 and vice-versa. Note that the previous equivalence holds even when any of
 243 both probabilities is scaled by a factor that does not depends on S , and
 244 consequently, the normalization constant of the LLHMM, $\mathcal{Z}_{\lambda}(O)$, defined
 245 in (29), can be removed from the right-hand side (34) without changing the
 246 equivalence. The proof of the equation is done in two steps by proving two
 247 implications: left to right, and right to left.

248 *4.1. From Generative to Discriminative Parameters*

249 Unlike the converse direction, it is relatively simple to prove that given a
 250 BHMM classifier for isolated word recognition, it can be reparameterized into
 251 a LLHMM. Recall that by definition of the LLHMM, if we set the log-linear
 252 parameters, λ , using the generative parameters, θ , as defined in (18)-(21),
 253 then we have that

$$p_{\theta}(O, S) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} \exp(\lambda' \mathbf{f}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})) = \mathcal{Z}_{\lambda}(O) p_{\lambda}(S | O). \quad (35)$$

254 Therefore, these two models when inserted into their corresponding classifiers
 255 in (34) produce proportional scores and, hence, select the same word or class.

256 *4.2. From Discriminative to Generative Parameters*

257 In this subsection, we prove the converse statement: that given a LLHMM
 258 classifier for isolated word recognition as defined in Sec. 3, an equivalent
 259 BHMM classifier exists. We begin expressing the right-hand model in (34)
 260 as follows

$$p_\lambda(S | O) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} \frac{h_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})}{\mathcal{Z}_\lambda(O)}, \quad (36)$$

261 where $h_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \exp(\lambda \mathbf{f}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}))$.

262 We start instantiating the feature $h_\lambda(\dots)$ in previous equations obtaining

$$h_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \exp(\lambda_S) \cdot h_\lambda(\mathbf{i}, \mathbf{q}; S) \cdot h_\lambda(O, \mathbf{k}; \mathbf{i}, \mathbf{q}, S) \quad (37)$$

263 with

$$h_\lambda(\mathbf{i}, \mathbf{q}; S) = \exp\left(\sum_{l,t} \lambda_{s_l q_t q_{t+1}}\right), \quad (38)$$

264 and

$$h_\lambda(O, \mathbf{k}; \mathbf{i}, \mathbf{q}, S) = \exp\left(\sum_{l,t} \lambda_{s_l q_t k_t} + \sum_{l,t,d} o_{td} \lambda_{s_l q_t k_t d}\right). \quad (39)$$

265 If we compare (9) expanded accordingly to (11) and (12), with (37) expanded
 266 with (38); then it is observed that each of the 3 terms in the right-hand side
 267 in (37) must be transformed, independently, into the corresponding term
 268 in (9).

269 Firstly, we transform $h_\lambda(O, \mathbf{k}; \mathbf{i}, \mathbf{q}, S)$ into $p_\theta(O, \mathbf{k} | \mathbf{i}, \mathbf{q}, S)$. Therefore,
 270 we need to transform the part of the discriminative parameters $\{\lambda_{c q k}\}$ and
 271 $\{\lambda_{c q k d}\}$ into the generative parameters $\{\tau; \mathbf{p}\}$, where τ is constrained as
 272 shown in (7). For doing that, (37) is multiplied and divided by $\exp(\sum_{l,t} \zeta_{s_l q_t})$,

273 and then, we rearrange the multiplication into (39) as follows

$$\exp\left(\sum_{l,t} \zeta_{s_l q_t}\right) h_\lambda(O, \mathbf{k}; \mathbf{i}, \mathbf{q}, S) = \exp\left(\sum_{l,t} (\lambda_{s_l q_t k_t} + \zeta_{s_l q_t}) + \sum_{l,t,d} o_{td} \lambda_{s_l q_t k_t d}\right), \quad (40)$$

274 whereas the division is moved into the second term in the right-hand side
275 of (37), yielding

$$\exp\left(-\sum_{l,t} \zeta_{s_l q_t}\right) h_\lambda(\mathbf{i}, \mathbf{q}; S) = \exp\left(\sum_{l,t} \bar{\lambda}_{s_l q_t q_{t+1}}\right), \quad (41)$$

276 where $\bar{\lambda}_{s q q'} = \lambda_{s q q'} - \zeta_{s q}$ will be used afterwards. The unknown parameters
277 $\{\zeta_{s q}\}$ are introduced to force the generative parameters $\{\tau_{c q}(k)\}$ to sum 1 in
278 the transformation.

From (16) and (40), and taking into account the constraints in (7), the solution must fulfill the following 3 constraints

$$\lambda_{c q k d} = \log \frac{p_{c q k d}}{1 - p_{c q k d}}, \quad (42)$$

$$\lambda_{c q k} + \zeta_{c q} = \log \tau_{c q}(k) + \xi_{c q}(k), \quad (43)$$

$$\sum_{k=1}^{K_{c q}} \tau_{c q}(k) = 1. \quad (44)$$

279 Then, from (42) we work out the value of $p_{c q k d}$

$$p_{c q k d} = \frac{\exp(\lambda_{c q k d})}{1 + \exp(\lambda_{c q k d})}, \quad (45)$$

280 and from (43) the value of $\tau_{c q}(k)$ is expressed as

$$\tau_{c q}(k) = \exp(\lambda_{c q k} - \xi_{c q}(k)) \exp(\zeta_{c q}), \quad (46)$$

281 where $\xi_{c q}(k)$ is defined as in (17) using the values of $\{p_{c q k d}\}$ defined in (45).

282 Although $\exp(\zeta_{c q})$ is still unknown, recall that it was introduced to tackle

283 the normalization constraint in (44), and then its value is worked out by
 284 replacing (46) in (44)

$$\exp(\zeta_{cq}) = \frac{1}{\sum_{k'=1}^{K_{cq}} \exp(\lambda_{cqq'} - \xi_{cq}(k'))}. \quad (47)$$

285 Finally, the exact value of $\tau_{cq}(k)$ is obtained by plugging (47) into (46)

$$\tau_{cq}(k) = \frac{\exp(\lambda_{cqq} - \xi_{cq}(k))}{\sum_{k'=1}^{K_{cq}} \exp(\lambda_{cqq'} - \xi_{cq}(k'))}. \quad (48)$$

286 Now we focus on transforming $\exp(-\sum_{l,t} \zeta_{s_l q_t}) h_\lambda(\mathbf{i}, \mathbf{q}; S)$ from (41) into
 287 $p_\theta(\mathbf{i}, \mathbf{q} | S)$ as defined in (11). This part of the proof is similar in conception to
 288 the proof given in Heigold et al. (2008b). Firstly we define a global transition
 289 matrix \mathcal{Q} as follows

$$\mathcal{Q}_{ij} = \begin{cases} \exp(\bar{\lambda}_{cqq'}) & i = f(c, q) \text{ and } j = f(c, q') \\ 1 & i = f(c, F) \\ 0 & \text{otherwise} \end{cases}, \quad (49)$$

290 where $f : \mathbb{N}^2 \mapsto \mathbb{N}$ is an injective function that maps each pair composed by
 291 a character and state, into a global index or state

$$f(c, q) = \begin{cases} B_c & q = I \\ B_c + q & 1 \leq q \leq M_c \\ B_c + M_c + 1 & q = F \end{cases}, \quad (50)$$

292 with M_c being the number of states for the symbol c , and $B_c = 1 + \sum_{n=1}^{c-1} (2 +$
 293 $M_n)$ being the number of preceding states to the first state of symbol c plus
 294 1.

295 Since all the values of \mathcal{Q} are not negative, accordingly to *Perron-Frobenius*
 296 *theorem*(Rao and Rao, 1998, p.473), the largest eigenvalue of \mathcal{Q} , ψ , is pos-
 297 itive and unique. Furthermore, the eigenvector associated to the largest
 298 eigenvalue, \mathbf{v} , has only positive coefficients, and obviously because of the
 299 eigenvector definition, \mathbf{v} satisfies

$$\sum_j \mathcal{Q}_{ij} v_j = \psi v_i, \quad \forall i = 1, \dots \quad (51)$$

300 Now, the transition generative parameters are defined as

$$a_c(q, q') = \frac{\mathcal{Q}_{f(c,q)f(c,q')} v_{f(c,q')}}{\psi v_{f(c,q)}} = \frac{\exp(\bar{\lambda}_{cqq'}) v_{f(c,q')}}{\psi v_{f(c,q)}}, \quad (52)$$

301 where $a_c(q, q')$ verifies the normalization constraint (7) because of (51). These
 302 parameters yield a probability proportional to that of (41) when used in (11)
 303 as the generative parameters of $p_\theta(\mathbf{i}, \mathbf{q} | S)$ (see Appendix A),

$$p_\theta(\mathbf{i}, \mathbf{q} | S) = \frac{1}{\psi^{T+L}} \left[\prod_l \frac{v_{f(s_l, F)}}{v_{f(s_l, I)}} \right] \frac{h_\lambda(\mathbf{i}, \mathbf{q}; S)}{\exp(\sum_{l,t} \zeta_{s_l q_t})}, \quad (53)$$

304 which is the equivalence we need but for the term $\frac{1}{\psi^{T+L}} \prod_l \frac{v_{f(s_l, F)}}{v_{f(s_l, I)}}$.

We can introduce this constant factor by multiplying and dividing (37)
 by it. The division is used in this part whereas the multiplication is added
 to the first term as follows

$$h_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \psi^T \cdot \exp(\bar{\lambda}_S) \cdot \frac{h_\lambda(\mathbf{i}, \mathbf{q}; S)}{\psi^{T+L} \exp(\sum_{l,t} \zeta_{s_l q_t})} \left[\prod_l \frac{v_{f(s_l, F)}}{v_{f(s_l, I)}} \right] \cdot \exp(\sum_{l,t} \zeta_{s_l q_t}) h_\lambda(O, \mathbf{k}; \mathbf{i}, \mathbf{q}, S) \quad (54)$$

305 with $\bar{\lambda}_S = \lambda_S + L \log \psi + \sum_{l=1}^L \log \frac{v_{f(s_l, I)}}{v_{f(s_l, F)}}$.

306 Finally, the last part of the proof consists in the transformation of $\exp(\bar{\lambda}_S)$
 307 into the word prior probabilities π_S . Similarly to the case of mixture coeffi-
 308 cients, we multiply and divide the numerator of (37) by an unknown constant,
 309 $\exp(\zeta)$. Since the constant $\exp(\zeta)$ is independent of the word S , it can be in-
 310 troduced into the right-hand side of (34). This constant is grouped together
 311 with $\exp(\bar{\lambda})$ as follows

$$\exp(\bar{\lambda}_S + \zeta) \tag{55}$$

Thus, taking into account (55) and the constraints (7), we have that following equalities must hold

$$\bar{\lambda}_S + \zeta = \log \pi_S \tag{56}$$

$$\sum_{S \in W} \pi_S = 1 \tag{57}$$

312 and the solution is found by following a similar procedure to that of the
 313 mixture coefficients

$$\pi_S = \frac{\exp(\bar{\lambda}_S)}{\sum_{S' \in W} \exp(\bar{\lambda}_{S'})}. \tag{58}$$

314 In summary, we have proven that for a given set of discriminative param-
 315 eters λ , a set of generative parameters can be defined, θ , by (45), (48), (52),

316 and (58); such that

$$\begin{aligned}
\arg \max_S p_\lambda(S | O) &= \arg \max_S \frac{\mathcal{Z}_\lambda(O)}{\exp(\zeta)} \exp(\zeta) p_\lambda(S | O) \\
&= \arg \max_S \exp(\zeta) \mathcal{Z}_\lambda(O) p_\lambda(S | O) \\
&= \arg \max_S \sum_{\mathbf{q}, \mathbf{i}, \mathbf{k}} \exp(\zeta) h_\lambda(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) \\
&= \arg \max_S \sum_{\mathbf{q}, \mathbf{i}, \mathbf{k}} \psi^T \exp(\bar{\lambda}_S + \zeta) \cdot \frac{h_\lambda(\mathbf{i}, \mathbf{q}; S)}{\psi^{T+L} \exp(\sum_{l,t} \zeta_{s_l q_t})} \left[\prod_{l=1}^L \frac{v_{f(s_l, F)}}{v_{f(s_l, I)}} \right] \\
&\quad \cdot \exp(\sum_{l,t} \zeta_{s_l q_t}) h_\lambda(O, \mathbf{k}; \mathbf{i}, \mathbf{q}, S) \\
&\Rightarrow \arg \max_S \sum_{\mathbf{q}, \mathbf{i}, \mathbf{k}} \psi^T \cdot \pi_S \cdot p_\theta(\mathbf{i}, \mathbf{q} | S) \cdot p_\theta(O, \mathbf{k} | \mathbf{i}, \mathbf{q}, S) \\
&= \arg \max_S \psi^T \sum_{\mathbf{q}, \mathbf{i}, \mathbf{k}} p_\theta(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}) = \arg \max_S p_\theta(O, S).
\end{aligned} \tag{59}$$

317 where by \Rightarrow we highlight the step of the proof that is not symmetric.

318 5. LLHMM Parameter Estimation

319 In contrast to generative models as BHMMs, in which parameter esti-
320 mation is usually carried out using the MLE criterion, there is not a unique
321 widespread criterion to find the optimal parameters for a class posterior
322 discriminative model such as the LLHMM proposed in this paper. Per-
323 haps, the most well-known criteria for discriminative parameter estimation
324 is the *maximum mutual information (MMI)*. Given a collection of samples
325 $\{(O_1, S_1), \dots, (O_N, S_N)\}$, the MMI criterion is defined as follows

$$F_{\text{MMI}}(\lambda) = \sum_{n=1}^N \log(p_\lambda(S_n | O_n)). \tag{60}$$

326 The optimal discriminative parameters, λ^* , are those that maximize F_{MMI} .

327 There are several algorithms for obtaining the parameters that maxi-
 328 mize (60) (Heigold et al., 2008a), but commonly the *Resilient back-propagation*
 329 (*RPROP*) *algorithm* (Riedmiller and Braun, 1993) is used (Giménez et al.,
 330 2011). The RPROP requires the computation of the gradient sign, for each
 331 parameter λ_m . The gradient of F_{MMI} is given by

$$\frac{\partial F_{\text{MMI}}(\lambda)}{\partial \lambda_m} = N_m(\lambda) - Q_m(\lambda) \quad (61)$$

332 where $N_m(\lambda)$ and $Q_m(\lambda)$ are expected counts defined as follows

$$N_m(\lambda) = \sum_{n=1}^N N_{nm}(\lambda), \quad Q_m(\lambda) = \sum_{n=1}^N Q_{nm}(\lambda), \quad (62)$$

333 with $N_{nm}(\lambda)$ and $Q_{nm}(\lambda)$ being the expected latent and class counts respec-
 334 tively. These counts are defined as follows

$$N_{nm}(\lambda) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} p_{\lambda}(\mathbf{i}, \mathbf{q}, \mathbf{k} \mid O_n, S_n) f_m(O_n, S_n, \mathbf{i}, \mathbf{q}, \mathbf{k}), \quad (63)$$

335 and

$$Q_{nm}(\lambda) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} \sum_S p_{\lambda}(S, \mathbf{i}, \mathbf{q}, \mathbf{k} \mid O_n) f_m(O_n, S, \mathbf{i}, \mathbf{q}, \mathbf{k}). \quad (64)$$

336 The probabilities $p_{\lambda}(\mathbf{i}, \mathbf{q}, \mathbf{k} \mid O, S)$ and $p_{\lambda}(S, \mathbf{i}, \mathbf{q}, \mathbf{k} \mid O)$ are computed as
 337 follows

$$p_{\lambda}(\mathbf{i}, \mathbf{q}, \mathbf{k} \mid O, S) = \frac{\exp(\lambda' \mathbf{f}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}))}{\mathcal{Z}_{\lambda}(O, S)}, \quad (65)$$

338 and

$$p_{\lambda}(S, \mathbf{i}, \mathbf{q}, \mathbf{k} \mid O) = p_{\lambda}(S \mid O) p_{\lambda}(\mathbf{i}, \mathbf{q}, \mathbf{k} \mid O, S) = \frac{\exp(\lambda' \mathbf{f}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k}))}{\mathcal{Z}_{\lambda}(O)}. \quad (66)$$

339 Finally, $\mathcal{Z}_\lambda(O)$ is the normalization constant for the model defined in (29)
 340 whereas $\mathcal{Z}_\lambda(O, S)$ is a joint normalization constant for the output and the
 341 word, which is likewise defined as follows

$$\mathcal{Z}_\lambda(O, S) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} \exp(\lambda' \mathbf{f}(O, S, \mathbf{i}, \mathbf{q}, \mathbf{k})). \quad (67)$$

342 The RPROP algorithm computes the sign of the gradient with the aid
 343 of these expected counts, and then, modifies the current parameters $\lambda^{(k)}$
 344 accordingly, so that a new estimate of the parameters is obtained, $\lambda^{(k+1)}$.
 345 The algorithm starts with a rough estimate of the parameters, $\lambda^{(0)}$, and it
 346 ends when either a maximum number of iterations have been reached, or the
 347 value of the objective function surpass a given threshold.

348 5.1. γ -MMI Criterion

349 A modification of the MMI criterion (60), the so-called γ -MMI criterion,
 350 leads to better performance (Schluter and Macherey, 1998; Povey, 2003). The
 351 γ -MMI is defined by introducing a scaling factor γ into the MMI criterion as
 352 follows

$$F_{\gamma\text{-MMI}}(\lambda) = \frac{1}{\gamma} \sum_{n=1}^N \log(p_{\lambda\gamma}(S_n | O_n)), \quad (68)$$

353 with $p_{\lambda\gamma}(S | O)$ defined as follows

$$p_{\lambda\gamma}(S | O) = \frac{[\mathcal{Z}_\lambda(O, S)]^\gamma}{\sum_R [\mathcal{Z}_\lambda(O, R)]^\gamma}. \quad (69)$$

354 The basic idea is to scale the likelihoods for each word in order to make
 355 the best words to compete one against the others even if the differences in
 356 probability are large. However, the reason why this idea outperforms the
 357 standard MMI is unclear. A possible hypothesis that we support is that it
 358 addresses some numerical problems related to the machine precision.

359 The gradient for the γ -MMI criterion in (68) is analogous to (61) but
 360 instead of using $Q_{nm}(\lambda)$, we now use $Q_{nm}^\gamma(\lambda)$ which is defined as follows

$$Q_{nm}^\gamma(\lambda) = \sum_{\mathbf{i}, \mathbf{q}, \mathbf{k}} \sum_S p_{\lambda\gamma}(S, \mathbf{i}, \mathbf{q}, \mathbf{k} | O_n) f_m(O_n, S, \mathbf{i}, \mathbf{q}, \mathbf{k}), \quad (70)$$

361 with the probability $p_{\lambda\gamma}(S, \mathbf{i}, \mathbf{q}, \mathbf{k} | O)$ defined as

$$p_{\lambda\gamma}(S, \mathbf{i}, \mathbf{q}, \mathbf{k} | O) = p_{\lambda\gamma}(S | O) p_\lambda(\mathbf{i}, \mathbf{q}, \mathbf{k} | S, O), \quad (71)$$

362 where the probabilities $p_{\lambda\gamma}(S | O)$ and $p_\lambda(\mathbf{i}, \mathbf{q}, \mathbf{k} | S, O)$ are defined in (69)
 363 and (65), respectively.

364 Fig. 2, summarizes the main idea behind the γ -MMI training criterion.
 365 It depicts the differences between the most probable word and the second
 366 most probable competitor for a LLHMM model (more details in Sec. 6). It
 367 is observed that these differences are of 44 points (in logarithmic scale) at
 368 the beginning, which corresponds to MLE. Additionally, the training sample
 369 is incorrectly classified at the first training iterations. Although, after 50
 370 iterations all the γ values correctly classify the sample; smaller values of γ
 371 induce larger difference between the correct class and its competitors.

372 5.2. Regularization

373 A common undesired property of all the proposed discriminative criteria
 374 is that they easily overfit the parameters. Even criteria specially designed
 375 to avoid outliers such as the power criterion suffer from overfitting. Since
 376 there is no clear way to smooth discriminatively trained models. A typical
 377 patch is to add a regularization term to the criterion itself

$$F_{C^*}(\lambda) = F_*(\lambda) - \frac{C}{2} \sum_m (\lambda_m^{(0)} - \lambda_m)^2, \quad (72)$$

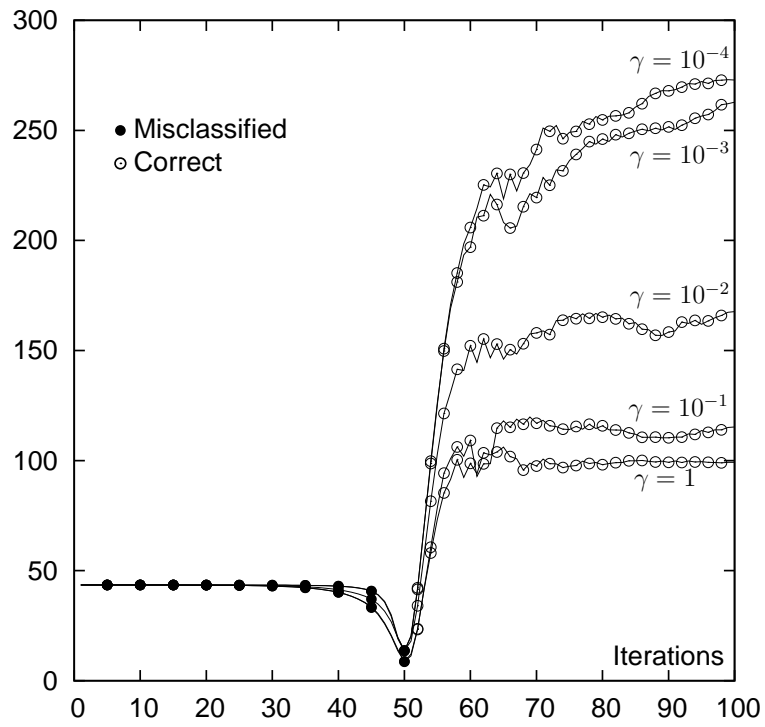


Figure 2: Differences (in logarithmic scale) between the most probable and the second most probable word for a given training sample (*cette*). Several values of the γ -MMI criterion are plotted. The most probable word changes at iteration 50 becoming the correct word.

378 with $F_*(\lambda)$ denoting the original criterion, namely F_{MMI} or F_γ ; and $\lambda^{(0)}$
 379 being either a reliable estimation of the parameters or simply $\mathbf{0}$.

380 The inclusion of the regularization term, only modifies the gradient in the
 381 following form

$$\frac{\partial F_{C^*}(\lambda)}{\partial \lambda_m} = \frac{\partial F_*(\lambda)}{\partial \lambda_m} + C(\lambda_m^{(0)} - \lambda_m) = N_m(\lambda) - Q_m(\lambda) + C(\lambda_m^{(0)} - \lambda_m), \quad (73)$$

382 where the expected counts, N_m and Q_m , are calculated as in the original
 383 criterion.

384 6. Experiments

385 In this section, we perform several experiments on the RIMES database
 386 of handwritten French letters (Grosicki et al., 2009), so that the performance
 387 of several discriminative training criteria for BHMM is assessed with respect
 388 to the generative training. Furthermore, we visually inspect several discrim-
 389 inative parameters by transforming them into their generative counterpart.

390 6.1. The RIMES Database

391 All the experiments were carried out on the protocol WR2 used in the
 392 handwritten word recognition competition of the ICDAR 2009. This protocol
 393 comprises 51 738 and 7 464 samples for training and testing, respectively. The
 394 lexicon used during the recognition comprises the words that occur in the test
 395 (1 612 words) and an alphabet of 81 characters. A three step preprocess was
 396 applied to all input images (Pastor i Gadea, 2007): gray level normalization,
 397 deslanting, and vertical size normalization.

398 Preprocessed images were first scaled in height to 30 pixels maintaining
 399 the aspect ratio, and then binarized with the Otsu’s method. A sliding

400 window of width 9 was applied centered on each column in order to extract a
401 sequence of 270-dimensional binary feature vectors. More precisely, for each
402 column the sliding window was horizontally centered, and then vertically
403 repositioned so that the center of the window is aligned with mass center of
404 the window before repositioning. Once realigned, the 9 binary vectors of the
405 window were concatenated in order to compose a binary feature vector of
406 dimension 270 which is fed into the BHMM or LLHMM as input.

407 *6.2. Experimental Setup*

408 In order to properly initialize the MMI training scheme the LLHMM
409 was initialized with a BHMM classifier trained with the EM (Baum-Welch)
410 algorithm (Rabiner and Juang, 1993), using the training scheme described
411 in Giménez et al. (2010). The best generative BHMM, which is composed
412 by $Q = 8$ states per character and $K = 64$ mixture components per state,
413 obtains an error of 21.2%.

414 Regarding to the discriminative training, the RPROP algorithm was used
415 for optimizing the criteria. The initial discriminative parameters were ob-
416 tained transforming the generative parameters of a BHMMs with $Q = 8$
417 states per character and $K = 26$ mixture components per state. Despite
418 the best generative results is obtained with $K = 64$ mixture components
419 per state, some works reported (Giménez et al., 2011) that the best classi-
420 fier obtained using MMI training has half (0.4 ratio) the number of mixture
421 component per state than its generative counterpart. Consequently, in pre-
422 liminary experiments we checked that the results obtained using the conven-
423 tional MMI criterion with $K = 26$ are similar or better to those obtained
424 increasing the value of K .

425 Finally, the proposed discriminative training criteria require to compute
 426 sums over all the words for calculating several values such as $\mathcal{Z}_\lambda(O)$ in (66).
 427 Consequently discriminative training algorithms become unfeasible in a straight
 428 implementation. For this reason we have approximated the sums over all the
 429 words by a beam pruning strategy together with a histogram pruning up to
 430 100 best hypothesis accordingly to $p(S | O)$.

431 6.3. Experiments

432 Firstly, we wanted to assess the repercussions of the regularization term
 433 in the conventional MMI criterion. For doing so, we scanned several values
 434 of the regularization parameter $C = \{0, 0.1, 1, 10, 100\}$ as introduced in (72),
 435 where $C = 0$ stands for not using the regularization at all. In Fig. 3, the
 436 classification error rate (CER) as a function of the number of RPROP iter-
 437 ations is plotted for different regularization values. In all cases, even with
 438 standard MMI, the CER decreases in a similar way, until iteration 60, where
 439 the best result is obtained. At this point, the behavior diverges depending
 440 on the precise value of C . If no regularization is applied ($C = 0$) the error
 441 becomes unstable and increases (overfits) as the training iterates. However,
 442 the larger the regularization parameter is, the less overtrained the model
 443 becomes, until that for $C = 10$ the error becomes stable while providing
 444 similar performance to that at iteration 60. Note that if the regularization
 445 parameter is further increased, a slight drop in performance is observed. As
 446 expected, the regularization term reduces the overfitting problem.

447 Fig. 3 also shows that the regularized MMI criterion obtains a CER
 448 around 19.5% using only $K = 26$ components per state. If we compare it
 449 with the best generative result, which is 21.2% and is obtained using $K = 64$,

450 we observe not only an improvement of 1.7 absolute points but also a reduc-
 451 tion on the number of parameters of 0.4, i.e. less than half the parameters
 are needed for such improvement.

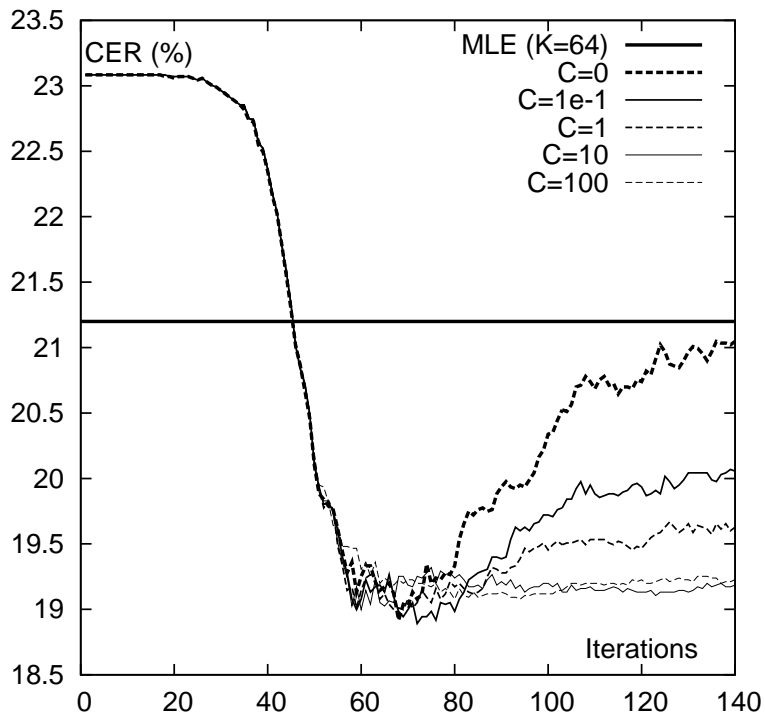


Figure 3: Classification error (in %) as a function of RPROP iterations for the regularized MMI criterion with several values of the regularization parameter C . Note that $C = 0$ stands for the non-regularized MMI.

452

453 For a deeper understanding of the MMI criterion, in Fig. 4 we depict the
 454 top 5 most probable words as a function of the training iterations (0,50,55,60,100)
 455 and the γ value of the γ -MMI criterion (1-MMII=MMI). We selected a com-
 456 mon example in which the MLE misclassifies the sample and the MMI learns
 457 how to discriminate it among the the other competitors. More precisely, for
 458 several training iterations the 5 most probable transcriptions are shown. In

459 addition, for each transcription the difference (in logarithmic scale) between
460 its score and the best score at that iteration is also shown. As expected,
461 the correct transcription (*cette*) gains relevance with the iterations, that is,
462 the training algorithm is modifying the model parameters in order to better
463 classify the sample. In particular, at the beginning there is a separation of
464 44 points between *cette* and the best transcription (*celle*). However, at some
465 point near to iteration 50 this situation is reverted, and from this point on
466 the score difference keeps increasing (see Fig. 2). A total of 60.3% of the
467 training samples that are misclassified by the MLE, are correctly classified
468 at the end of the last MMI iteration. In contrast, only 1.3% of the correctly
469 classified samples by the *MLE* are misclassified at the end of the training
470 process.

471 In Fig. 6, we explored several values of γ , ranging from standard MMI
472 ($\gamma = 2$) to 10^{-4} , using the best regularization term obtained in the previous
473 experiment $C = 10$. In the previous experiment the 100-best words were
474 recalculated every 10 iterations, however, with the γ -MMI we observed a
475 severe overfitting. Consequently, we repeated the experiments recalculating
476 the best words every iteration. It is observed in Fig. 6 that the modified
477 γ -MMI obtains a very competitive performance in terms of CER (15%) if
478 applied properly. If we compare the best result in Fig. 6 with the best
479 generative result, the former obtains an improvement boost of more than 6
480 absolute points with respect to the latter. In Fig. 4 the behavior of the γ -MMI
481 can be checked for a training sample. As we can see, the use of small values
482 of γ leads to an increase of the separation between classes, which is consistent
483 with the idea that the γ is increasing the competition between classes during

Original image



Preprocessed without sliding window



	MLE		γ -MMI ($\gamma = 1$)							
Iterations	1		50		55		60		100	
1-best	-	celle	-	celle	-	cette	-	cette	-	cette
2-best	44	cette	6	cette	73	celle	72	celle	100	celle
3-best	447	Cette	376	dette	342	dette	359	dette	354	dette
4-best	467	dette	406	Cette	382	Cette	388	Cette	393	Cette
5-best	499	celles	497	celles	564	celles	542	celles	485	geste

	MLE		γ -MMI ($\gamma = 10^{-1}$)							
Iterations	1		50		55		60		100	
1-best	-	celle	-	celle	-	cette	-	cette	-	cette
2-best	44	cette	7	cette	79	celle	94	celle	116	celle
3-best	447	Cette	375	dette	340	dette	358	dette	336	dette
4-best	467	dette	409	Cette	385	Cette	403	Cette	393	Cette
5-best	499	celles	497	celles	570	celles	556	geste	459	geste

	MLE		γ -MMI ($\gamma = 10^{-3}$)							
Iterations	1		50		55		60		100	
1-best	-	celle	-	celle	-	cette	-	cette	-	cette
2-best	44	cette	6	cette	128	celle	207	celle	260	celle
3-best	447	Cette	403	dette	413	dette	417	dette	440	dette
4-best	467	dette	516	celles	663	Cette	631	Cette	657	Cette
5-best	499	celles	530	Cette	667	celles	784	celles	828	telle

Figure 4: γ -MMI behavior on a training sample for several values of γ . The figures stand for the difference (in logarithmic scale) between each n -best word and the best transcription at each iteration. Bold words highlight the position of the correct word *cette*.

484 the training process. For example, at iteration 100 the separation between
485 the two best hypothesis using $\gamma = 1$ is 100 points, while using $\gamma = 10^{-3}$ the
486 separation increases up to 260 points.

487 Fig. 5 depicts a similar experimentation to that of Fig. 3 but for several
488 *test samples*. The first sample (*vous*) is a sample that is misclassified by the
489 MLE model and it is correctly classified using γ -MMI criterion. The remain-
490 ing two samples are correctly classified by the MLE criterion. However, the
491 first one is finally misclassified by the discriminative model, while the second
492 one remains correctly classified. It is worth noting, that these three cases
493 represent the 10.2%, 2.2% and 74.7% of the test set, respectively.

494 As discussed before, all experiments were carried out using $K = 26$ com-
495 ponents per state. In order to better compare the performance of the MLE
496 and γ -MMI criteria we carried out a final experiment, in which both criteria
497 are tested using several components per state $K \in \{1, 4, 16, 64\}$. For the
498 γ -MMI criteria the best parameters from previous experiments were used
499 ($\gamma = 10^{-3}$ and $C = 10$). Results are shown in Fig. 7.

500 From the results reported in Fig. 7 it is clear that γ -MMI outperforms
501 MLE in all cases. The improvement of the MMI decreases as the number of
502 components increases. For example, the improvement using $K = 1$ is about
503 20 points while using $K = 64$ is about 5 points. It is worth noting, that the
504 best result in this figure is 15.2% which is achieved using $K = 16$ components
505 and it is very similar to the best result obtained with $K = 26$, which we chose
506 for all the previous experimentation.

507 Finally, a visual inspection of some Bernoulli prototypes for several train-
508 ing criteria is given in Fig. 8. The Bernoulli prototypes for letters *e* and *s* are

Incorrect \rightarrow Correct

Iterations	MLE		γ -MMI ($\gamma = 10^{-3}$)							
	0		40	45	50	100				
1-best	-	virus	-	virus	-	virus	-	vous	-	vous
2-best	26	Vous	35	vous	15	vous	16	virus	89	virus
3-best	44	vous	40	Vous	50	Vous	93	Vous	318	Vous
4-best	82	bruits	118	bruits	165	bruits	268	bruits	350	nous
5-best	231	plus	243	plus	261	plus	323	plus	419	viens

Correct \rightarrow Incorrect

Iterations	MLE		γ -MMI ($\gamma = 10^{-3}$)							
	0		40	45	50	100				
1-best	-	Suite	-	Suite	-	Suite	-	suite	-	suite
2-best	97	suite	77	suite	50	suite	16	Suite	247	Suite
3-best	357	Socit	376	Socit	395	seule	336	seule	465	seule
4-best	405	socit	413	socit	402	Socit	457	socit	698	sant
5-best	428	Sant	431	seule	425	socit	470	Socit	702	suis

Correct \rightarrow Correct

Iterations	MLE		γ -MMI ($\gamma = 10^{-3}$)							
	0		40	45	50	100				
1-best	-	que	-	que	-	que	-	que	-	que
2-best	239	due	233	due	221	due	188	due	167	due
3-best	350	dire	371	dire	396	dire	438	dire	753	dire
4-best	539	grise	590	grise	620	date	628	date	810	date
5-best	595	avez	611	date	639	avez	659	d'une	930	quel

Figure 5: γ -MMI behavior on three selected *test* samples. The figures stand for the difference (in logarithmic scale) between each n -best word and the best transcription at each iteration. Bold words highlight the position of the correct word *cette*.

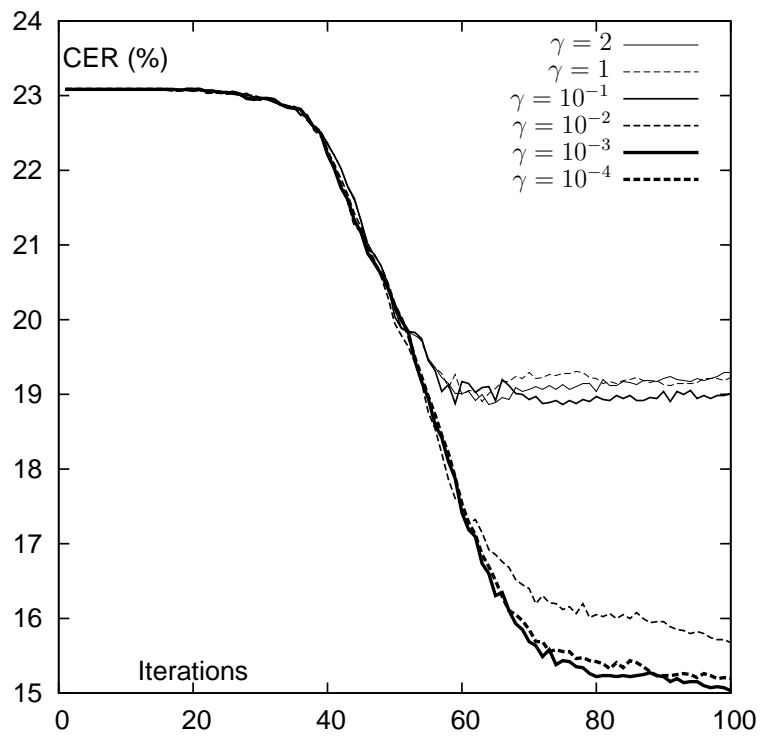


Figure 6: Classification error (in %) as a function of RPROP iterations for the modified γ -MMI criterion with regularization $C = 10$ and several values of γ . Note that $\gamma = 1$ corresponds to the standard MMI criterion.

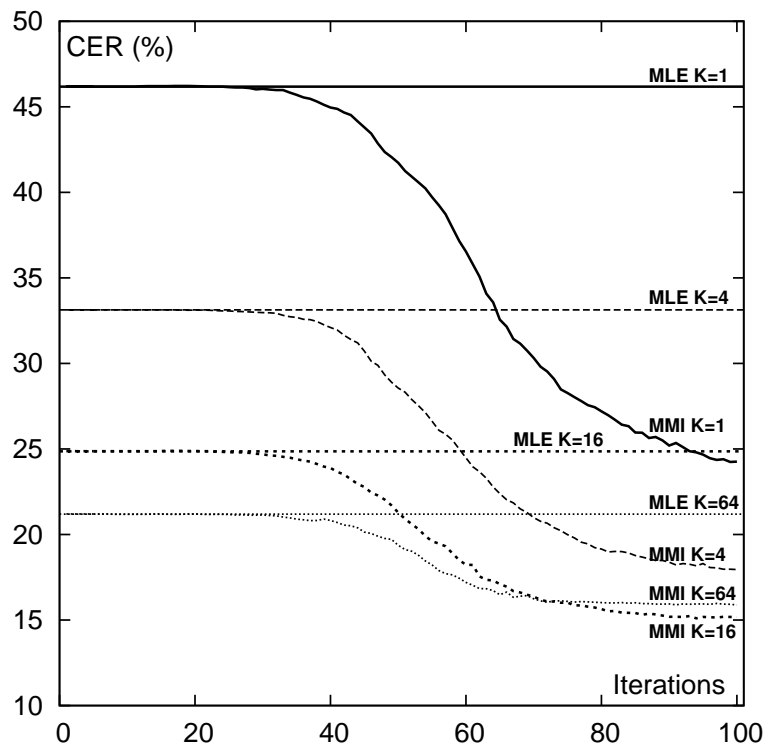


Figure 7: MLE and γ -MMI ($\gamma = 10^{-3}C = 10$) criteria comparison using several components per state.

509 shown, where the columns represent states, and rows represent mixture com-
510 ponents. Provided that the number of mixture components in each state is
511 large ($K = 26$) we have selected the 4 components with the highest mixture
512 coefficients when trained using the MLE criterion. Prototypes are plotted
513 for 3 different training criteria (from left to right): MLE training; the γ -MMI
514 with $\gamma = 10^{-3}$ and regularization $C = 10$; and the conventional MMI train-
515 ing without regularization. It is worth noting that the MLE prototypes are
516 the initial prototypes for both represented discriminative training criteria.

517 It is observed that the prototypes without regularization are apparently
518 a noise version of the MLE prototypes, however we know that they have
519 a better performance when classifying. A further observation reveals that
520 discriminative training is focused on modifying those pixels that discriminate
521 the most while keeping the remaining pixels unmodified. These unmodified
522 pixels are those that keep the same state (0 or 1) for many words. When no
523 regularization is employed, a pixel that discriminates a single training sample
524 can be set to 1, however, those spurious pixels are eliminated by adding the
525 regularization term.

526 7. Concluding Remarks and Future Work

527 In this work, we presented a log-linear HMM (LLHMM) to recognize
528 isolated handwritten words that directly deals with binarized images with-
529 out the need of a sophisticated feature extraction process. This model has
530 been proved to be equivalent to Bernoulli HMMs (BHMMs), and in this way,
531 we have provided a framework for discriminatively training BHMMs. Fur-
532 thermore, this allows us to visually inspect and understand discriminative

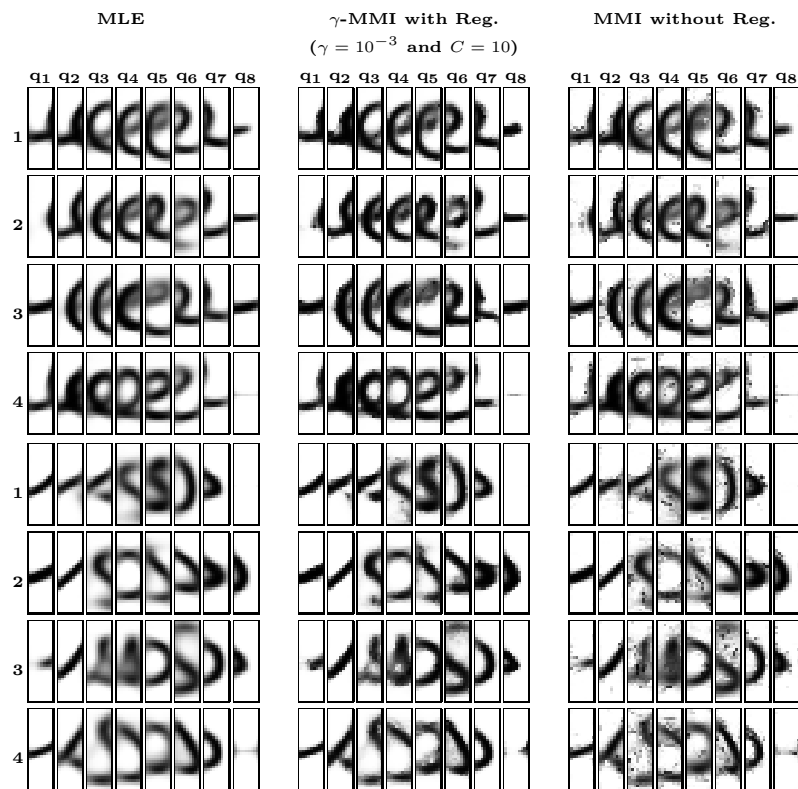


Figure 8: Bernoulli prototypes of letters e and s using three different training criteria (from left to right): MLE, γ -MMI with regularization and MMI without regularization.

533 parameters by transforming them into generative ones.

534 Two discriminative training criteria have also been analyzed for the LLHMM
535 model: conventional MMI and γ -MMI. We tried all of them discussing prob-
536 lems (over-fitting, computational cost) and some typical approximation to
537 those problems (regularization term, pruning techniques). All these methods
538 have been tested over the well-known RIMES database of handwritten French
539 words. Furthermore, in all cases discriminative training clearly outperformed
540 the conventional MLE training. In particular, very competitive results were
541 obtained using the γ -MMI training scheme which obtained nearly 15% of
542 CER, or in other words an improvement of more than 6% of absolute points
543 with respect to the generative counterpart. However, there are many more
544 discriminative training criteria such as margin-based or minimum phoneme
545 error. As future work we plan to implement and adapt these discriminative
546 criteria to the proposed model.

547 The best result obtained in this work on the considered task of the RIMES
548 database is 15%, which to our knowledge is the best result reported using
549 HMMs and without system combination. If we compare our system with
550 the results of the ICDAR 2009 (Grosicki and El Abed, 2009), our system
551 would be positioned in the third position and very close to the second sys-
552 tem (13.9%), which is in fact a combination of hybrid HMM/MLP systems,
553 and far from the first system (6.8), which is a system based on a hierarchy of
554 multidimensional recurrent neural networks, and has shown to be extremely
555 competitive in this task. Moreover, if we compare our result with the results
556 reported on the same task in Bianne-Bernard et al. (2011), it is observed that
557 our system outperforms the results of the three systems presented on that

558 paper: a dynamic context-independent system based on HMMs (24.5%), a
559 dynamic context-dependent system based on HMMs (19.6%), and a hybrid
560 HMM/neural network system (20.5%). However, when the three systems are
561 combined an error of 10.9% is obtained. Consequently, as future work we
562 plan to combine the proposed discriminative BHMMs system, with the con-
563 ventional generative BHMMs system and other state of the art systems, as
564 for instance those based on recurrent neural networks (Graves and Schmid-
565 huber), in order to measure the impact of discriminative BHMMs when com-
566 bined with other systems.

567 Finally, we intend to extend all the work developed in this paper to con-
568 tinuous HTR, that is, a discriminative BHMM in which the words are re-
569 placed by word sequences, and hence, the prior probabilities are replaced by
570 a language model.

571 **Acknowledgment**

572 Work supported by the EC (FEDER/FSE) and the Spanish MEC/MICINN
573 under the MIPRCV “Consolider Ingenio 2010” program (CSD2007-00018),
574 iTrans2 (TIN2009-14511) and MITTRAL (TIN2009-14633-C03-01) projects.
575 Also supported by the IST Programme of the European Community, under
576 the PASCAL2 Network of Excellence, IST-2007-216886, and by the Spanish
577 MITyC under the erudito.com (TSI-020110-2009-439).

578 **Appendix A. Discriminative to generative transition probabilities**

579 In this appendix, we prove that the parameters in (52) yield a probability
580 proportional to that of (41) when used in (11) as the generative parameters

581 of $p_\theta(\mathbf{q}, \mathbf{i} \mid S)$. In order to clarify this, we plug the parameters as computed
 582 in (52) into (11) yielding

$$\prod_{l=1}^L \left[\frac{\exp(\bar{\lambda}_{s_l I q_{i_l}}) v_{f(s_l, q_{i_l})}}{\psi v_{f(s_l, I)}} \cdot \prod_{t=i_l}^{i_{l+1}-2} \frac{\exp(\bar{\lambda}_{s_l q_t q_{t+1}}) v_{f(s_l, q_{t+1})}}{\psi v_{f(s_l, q_t)}} \cdot \frac{\exp(\bar{\lambda}_{s_l q_{i_{l+1}-1} F}) v_{f(s_l, F)}}{\psi v_{f(s_l, q_{i_{l+1}-1})}} \right], \quad (\text{A.1})$$

583 where by grouping elements we get

$$\frac{1}{\psi^{T+L}} \frac{h_\lambda(\mathbf{i}, \mathbf{q}; S)}{\exp(\sum_{l,t} \zeta_{s_l q_t})} \prod_{l=1}^L \left[\frac{v_{f(s_l, F)}}{v_{f(s_l, I)}} \frac{v_{f(s_l, q_{i_l})}}{v_{f(s_l, q_{i_{l+1}-1})}} \prod_{t=i_l}^{i_{l+1}-2} \frac{v_{f(s_l, q_{t+1})}}{v_{f(s_l, q_t)}} \right]. \quad (\text{A.2})$$

584 Note that, in each segment l the telescope product over $\frac{v_{j'}}{v_j}$ is equal to $\frac{v_{f(s_l, F)}}{v_{f(s_l, I)}} \cdot 1$,
 585 and then equation (A.2) is reduced to

$$p_\theta(\mathbf{i}, \mathbf{q} \mid S) = \frac{1}{\psi^{T+L}} \left[\prod_{l=1}^L \frac{v_{f(s_l, F)}}{v_{f(s_l, I)}} \right] \frac{h_\lambda(\mathbf{i}, \mathbf{q}; S)}{\exp(\sum_{l,t} \zeta_{s_l q_t})}. \quad (\text{A.3})$$

586 References

- 587 Bianne-Bernard, A.L., Menasri, F., Al-Hajj Mohamad, R., Mokbel, C., Ker-
 588 morvant, C., Likforman-Sulem, L., 2011. Dynamic and contextual informa-
 589 tion in hmm modeling for handwritten word recognition. *Pattern Analysis*
 590 *and Machine Intelligence*, IEEE Transactions on 33, 2066 –2080.
- 591 Dempster, A.P., Laird, N.M., Rubin, D.B., 1977. Maximum Likelihood from
 592 Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical*
 593 *Society. Series B (Methodological)* 39, 1–38.
- 594 Duda, R.O., Hart, P.E., 1973. *Pattern Classification and Scene Analysis*. J.
 595 Wiley and Sons.

- 596 Giménez, A., Alkhoury, I., Juan, A., 2010. Windowed Bernoulli Mixture
597 HMMs for Arabic Handwritten Word Recognition, in: ICFHR' 10, Kolkata
598 (India). pp. 533–538.
- 599 Giménez, A., Andrés-Ferrer, J., Juan, A., Serrano, N., 2011. Discriminative
600 Bernoulli Mixture Models for Handwritten Digit Recognition, in: ICDAR'
601 11, Beijing (China). pp. 558–562.
- 602 Giménez, A., Juan, A., 2009. Embedded Bernoulli Mixture HMMs for Hand-
603 written Word Recognition, in: ICDAR' 09, Barcelona (Spain). pp. 896–900.
- 604 Graves, A., Schmidhuber, J., . Offline Handwriting Recognition with Multi-
605 dimensional Recurrent Neural Networks.
- 606 Grosicki, E., Carree, M., Brodin, J.M., Geoffrois, E., 2009. Results of the
607 RIMES Evaluation Campaign for Handwritten Mail Processing, in: IC-
608 DAR' 09, Barcelona(Spain). pp. 941–945.
- 609 Grosicki, E., El Abed, H., 2009. ICDAR 2009 Handwriting Recognition
610 Competition, in: ICDAR '09, Barcelona (Spain). pp. 1398–1402.
- 611 Heigold, G., Deselaers, T., Schlüter, R., Ney, H., 2008a. A GIS-like training
612 algorithm for log-linear models with hidden variables, in: ICASSP' 08, Las
613 Vegas (USA). pp. 4045–4048.
- 614 Heigold, G., Lehnen, P., Schluter, R., Ney, H., 2008b. On the equivalence of
615 Gaussian and log-linear HMMs, in: INTERSPEECH' 08, Brisbane (Aus-
616 tralia). pp. 273–276.

- 617 Märgner, V., El Abed, H., 2010. ICFHR 2010 - Arabic Handwriting Recog-
618 nition Competition, in: ICFHR' 10, Kolkata (India). pp. 709–714.
- 619 Pastor i Gadea, M., 2007. Aportaciones al reconocimiento automático de
620 texto manuscrito. Ph.D. thesis. Dep. de Sistemes Informàtics i Com-
621 putació. València, Spain.
- 622 Povey, D., 2003. Discriminative Training for Large Vocabulary Speech Recog-
623 nition. Ph.D. thesis. Cambridge University Engineering Dept.
- 624 Rabiner, L., Juang, B.H., 1993. Fundamentals of speech recognition.
625 Prentice-Hall.
- 626 Rao, C.R., Rao, M.B., 1998. Matrix algebra and its applications to statistics
627 and econometrics. World Scientific.
- 628 Riedmiller, M., Braun, H., 1993. A Direct Adaptive Method for Faster Back-
629 propagation Learning: The RPROP Algorithm, in: IEEE International
630 Conference on Neural Networks, pp. 586–591.
- 631 Schluter, R., Macherey, W., 1998. Comparison of discriminative training
632 criteria, in: ICASSP' 98, pp. 493–496 vol.1.