



ELSEVIER

Available online at www.sciencedirect.com

Pattern Recognition Letters xxx (2007) xxx–xxx

Pattern Recognition
Letterswww.elsevier.com/locate/patrec

A new look at discriminative training for hidden Markov models

Xiaodong He *, Li Deng *

Microsoft Research, 1 Microsoft Way, Redmond, WA 98052, United States

6 Abstract

Discriminative training for hidden Markov models (HMMs) has been a central theme in speech recognition research for many years. One most popular technique is minimum classification error (MCE) training, with the objective function closely related to the empirical error rate and with the optimization method based traditionally on gradient descent. In this paper, we provide a new look at the MCE technique in two ways. First, we develop a non-trivial framework in which the MCE objective function is re-formulated as a rational function for multiple sentence-level training tokens. Second, using this novel re-formulation, we develop a new optimization method for discriminatively estimating HMM parameters based on growth transformation or extended Baum–Welch algorithm. Technical details are given for the use of lattices as a rich representation of competing candidates for the MCE training.

© 2007 Elsevier B.V. All rights reserved.

Keywords: Hidden Markov model; Discriminative learning; Minimum classification error; Extended Baum–Welch algorithm; Growth transformation

17 1. Introduction

Hidden Markov models (HMMs) have been a well established framework for a variety of pattern recognition applications, including, most prominently, speech recognition applications (Rabiner and Juang, 1993; Bahl et al., 1987; Deng and O’Shaughnessy, 2003). One most attractive feature of the HMM framework is that its parameters can be learned automatically from the training data. In early days of HMMs, the parameters were learned by the maximum likelihood (ML) criterion based on the EM algorithm (e.g., Bahl et al., 1987; Rabiner and Juang, 1993). Improvement of parameter learning beyond ML has been pursued for many years (Brown, 1987; Chou, 2003; Deng et al., 2005; Deng et al., 2005; Gopalakrishnan et al., 1991; He and Chou, 2003; Juang and Katagiri, 1992; Juang et al., 1997; Macherey et al., 2005; McDermott et al., in press; Normandin, 1991; Povey and Woodland, 2002; Povey et al., 2003; Povey, 2004; Povey et al., 2004; Rathinavalu

and Deng, 1998; Schluter et al., 2001), based on the concept of discrimination against classes, in contrast to maximizing likelihood of each individual class. The reason behind discriminative training is that complete knowledge of speech data distributions is lacking and training data is always limited. It is not until recently that discriminative training has shown uniform success in speech recognition over virtually all tasks, including especially large tasks (e.g., Woodland and Povey, 2000; Povey, 2004).

Among several types of discriminative training for HMMs, one prominent type is minimum classification error (MCE) training (Chou, 2003; Juang and Katagiri, 1992; Juang et al., 1997; He and Chou, 2003; Macherey et al., 2005; McDermott et al., in press; Roux and McDermott, 2005; Rathinavalu and Deng, 1998). The essence of MCE is to define the objective function for optimization that is closely related to the empirical classification errors. This is more desirable than other types of discriminative training that are less closely related to the classification errors. The conventional MCE has been based on the sequential gradient-descent based technique, named generalized probabilistic descent (GPD), which optimizes the objective function as a highly complex function of the HMM parameters.

* Corresponding authors. Tel.: +1 425 706 4939 (Xiaodong He); +1 425 706 2719 (Li Deng); fax: +1 425 936 7329.

E-mail addresses: xiaoh@microsoft.com (X. He), deng@microsoft.com (L. Deng).

Another significant advance in discriminative training is the development and application of a special type of optimization technique, called growth transformation (GT) or extended Baum–Welch (EBW) algorithm when it is used for HMM parameter estimation. GT is an iterative optimization scheme where if the parameter set A is subject to a transformation $A = T(A')$, then the objective function “grows” in its value $O(A) > O(A')$ unless $A = A'$. In (Gopalakrishnan et al., 1991), GT/EBW was developed for rational functions such as the mutual information as the optimization criterion. Maximization of mutual information (MMI) as a form of discriminative criterion for the discrete HMM was described in (Gopalakrishnan et al., 1991). This has been extended to the continuous-density HMM in (Normandin, 1991; Gunawardana and Byrne, 2001). The significance of GT/EBW lies in its effectiveness and closed-form parameter updating for large-scale optimization problems with difficult objective functions. Compared with the gradient based techniques which often require special and delicate care for tuning the parameter-dependent learning rate, GT/EBW avoids such requirements and with the closed-form updating formula it is generally faster in reaching algorithm convergence.

Mutual information is naturally in the form of a rational function and MMI is obviously suited to GT/EBW optimization. However, as a discriminative criterion, it is only indirectly related to classification errors. On the other hand, MCE as a discriminative criterion is closely related to classification errors, but it is not naturally in the form of a rational function when there are multiple utterance tokens in the training data. Hence, it has been a tradition to use the gradient-descent techniques (GPD) for optimizing the MCE criterion (Chou, 2003; Juang et al., 1997; McDermott et al., in press; Rathinavalu and Deng, 1998). In this paper, we break this long-held tradition and take a fresh look at the MCE. This new analysis and formulation of the MCE covers two main issues. First we re-examine the MCE criterion. Second the results of the re-examination permit the use of the new GT/EBW optimization technique for optimizing the MCE criterion with respect to the HMM parameters.

The organization of this paper is as follows. In Section 2, an overview of the traditional MCE is provided. Then, in Section 3, we reformulate the MCE criterion (with multiple training tokens) into a rational functional form. We provide a rigorous proof by induction for the correctness of the rational functional form. Given this non-trivial reformulation, in Section 4, we present in detail a novel GT/EBW based optimization technique for estimating the parameters of the Gaussian HMMs. In Section 5, the lattice-based MCE training is described, and a summary is given in Section 6.

2. Overview of minimum classification error (MCE) training

We denote by A the parameter set of the generative model expressed in terms of a joint statistical distribution

$$p_A(X, S) = p_A(X|S)P(S), \quad (1) \quad 115$$

on the observation training data sequence X and on the corresponding label sequence S , where we assume the parameters in the “language model” $P(S)$ are not subject to optimization. We use $r = 1, \dots, R$ as the index for “token” (e.g., a single sentence or utterance) in the training data, and each token consists of a “string” of an observation data sequence: $X_r = x_{r,1}, \dots, x_{r,T_r}$, with the corresponding label (e.g., word) sequence: $S_r = w_{r,1}, \dots, w_{r,N_r}$. That is, S_r denotes correct label sequence for token r . Further, we use s_r to denote all possible label sequences for the r th token, including the correct label sequence S_r and all other incorrect label sequences.

MCE learning was originally introduced for multiple-category classification problems where the smoothed error rate is minimized for isolated “tokens” (Juang and Katagiri, 1992). It was later generalized to minimize the smoothed “sentence token” or string-level error rate (Juang et al., 1997; Chou, 2003), which is known as “embedded MCE”.

The MCE objective function is defined first based on a set of class discriminant functions and a special type of loss function. Then the model is estimated to minimize the expected loss that is closely related to the recognition error rate of the classifier.

In embedded MCE training, a set of discriminant functions is first defined based on the correct string S_r and the N most confusable competing strings, $s_{r,1}, \dots, s_{r,N}$. Define the top N best competing strings as

$$s_{r,1} = \arg \max_{s_r: s_r \neq S_r} \{\log p_A(X_r, s_r)\},$$

$$s_{r,i} = \arg \max_{s_r: s_r \neq S_r, s_r \neq s_{r,1}, \dots, s_{r,i-1}} \{\log p_A(X_r, s_r)\} \quad i = 2, \dots, N. \quad 145$$

Then, the discriminant functions for the correct string and the N competing strings take the form of

$$g_{s_r}(X_r; A) = \log p_A(X_r, s_r), \quad s_r \in \{S_r, s_{r,1}, \dots, s_{r,N}\}. \quad 149$$

And the decision rule for the recognizer or classifier is the one that for the observation data sequence, X_r ,

$$C(X_r) = s_r^* \quad \text{if } s_r^* = \arg \max_{s_r} g_{s_r}(X_r; A). \quad 153$$

Next, a misclassification measure in MCE is defined. For the general N -best MCE training, the following misclassification measure has been widely used (Juang et al., 1997):

$$d_r(X_r, A) = -\log p_A(X_r, S_r)$$

$$+ \log \left\{ \frac{1}{N} \sum_{s_r: s_r \neq S_r} \exp[\eta \log p_A(X_r, s_r)] \right\}^{\frac{1}{\eta}}. \quad (2) \quad 160$$

This misclassification measure function emulates the decision rule, i.e., $d_r(X_r, A) \geq 0$ implies misclassification and $d_r(X_r, A) < 0$ implies a correct classification. The second term in (2) is a soft-max function, which counts the scores of all N competitive candidates. It can be looked as an average over the scores of competitive candidates

167 weighted based on their individual significance. Moreover,
168 this misclassification measure can be closely approximated
169 by the following simpler form:

$$170 \quad d_r(X_r, A) = -\log p_A(X_r, S_r) + \log \sum_{s_r, s_r \neq S_r} w(s_r) \cdot p_A(X_r, s_r),$$

172 (3)

173 where $w(s_r)$ is a non-negative weighting factor for compet-
174 itive string s_r . Note that the sum of $w(s_r)$ is not necessarily
175 equal to one.

176 Finally, to define the objective function of MCE, a loss
177 function for a single sentence token or string X_r is estab-
178 lished, as originally proposed in (Juang and Katagiri,
179 1992; Juang et al., 1997), in the following form:

$$181 \quad l_r(d_r(X_r, A)) = \frac{1}{1 + e^{-\alpha d_r(X_r, A) + \beta}} = \frac{1}{1 + e^{-d_r(X_r, A)}},$$

182 (4)

182 where we assume $\alpha = 1$, $\beta = 0$ for simplicity in exposition
183 without loss of generality. This loss function emulates the
184 zero-one recognition error count function, i.e., when
185 $d_r(X_r, A)$ is larger than zero, which implies an incorrect rec-
186 ognition, the loss function approaches to one, which essen-
187 tially becomes a recognition error count.

188 With the misclassification measure in the form of (3), the
189 loss function for the N -best version of MCE becomes:

$$192 \quad l_r(d_r(X_r, A)) = \frac{\sum_{s_r, s_r \neq S_r} w(s_r) p_A(X_r, s_r)}{\sum_{s_r, s_r \neq S_r} w(s_r) p_A(X_r, s_r) + p_A(X_r, S_r)}$$

$$= \frac{\sum_{s_r, s_r \neq S_r} w(s_r) p_A(X_r, s_r)}{\sum_{s_r} w(s_r) p_A(X_r, s_r)}.$$

193 (5)

193 The last step is obtained after the assignment of
194 $w(S_r) \equiv 1$ for the correct string S_r .

195 Given the loss function for one sentence token r in (5),
196 the empirical loss function over the whole training set with
197 all R training tokens becomes:

$$200 \quad L(A) = \sum_{r=1}^R l_r(d_r(X_r, A)).$$

201 (6)

201 Therefore, (6) is closely related to the empirical recogni-
202 tion error rate and is the objective function to minimize in
203 MCE. The traditional MCE methods minimize the loss
204 function via the technique of probabilistic gradient descent
205 or GPD, which we refer the readers to an excellent review
206 in (Chou, 2003).

207 3. A new look at MCE – optimization criterion

208 We now take a new look at MCE in terms of its optimi-
209 zation criterion as expressed in (6). Minimizing the overall
210 loss function of $L(A)$ in (6) is to the same as maximizing the
211 following equivalent objective function:

212

$$O(A) = R - L(A) = \sum_{r=1}^R \left[1 - \frac{\sum_{s_r, s_r \neq S_r} w(s_r) p_A(X_r, s_r)}{\sum_{s_r} w(s_r) p_A(X_r, s_r)} \right]$$

$$= \sum_{r=1}^R \frac{w(S_r) p_A(X_r, S_r)}{\sum_{s_r} w(s_r) p_A(X_r, s_r)}.$$

214 (7)

215 Importantly, (7) is a sum of rational functions rather than
216 a rational function in itself, and hence it would not be
217 directly amenable to GT/EBW for its optimization. The dif-
218 ficulty of formulating a rational function and the desire of
219 moving away from gradient descent have been discussed in
220 (Povey, 2004). In this section, we directly tackle this diffi-
221 culty and re-formulate the MCE objective function of (7)
222 as a true rational function of the following specific form:

$$223 \quad O(A) = \frac{\sum_{s_1 \dots s_R} w(s_1 \dots s_R) p_A(X_1 \dots X_R, s_1 \dots s_R) C(s_1 \dots s_R)}{\sum_{s_1 \dots s_R} w(s_1 \dots s_R) p_A(X_1 \dots X_R, s_1 \dots s_R)},$$

224 (8) 225

226 where $w(s_1 \dots s_R) = \prod_{r=1}^R w(s_r)$ and $C(s_1 \dots s_R) = \sum_{r=1}^R C(s_r)$,
227 $C(s_r) = \delta(s_r, S_r)$. Here, $\delta(s_r, S_r)$ is the Kronecker delta func-
228 tion that equals one if $s_r = S_r$, and zero otherwise. Note that
229 $w(s_1, \dots, s_R)$ and $C(s_1, \dots, s_R)$ are quantities not relevant to
230 A . In (8), $X = X_1, \dots, X_R$ denotes the collection of all obser-
231 vation data sequences in all R training tokens, and
232 $p_A(X_1, \dots, X_R, s_1, \dots, s_R)$ is the joint distribution for all train-
233 ing data and their corresponding label sequence assign-
234 ments $s = s_1, \dots, s_R$.

235 We now provide a rigorous proof that (7) and (8) are
236 equivalent. We use the induction method for the proof in
237 the following two steps.

238 (1) We prove the equivalence of (7) and (8) when there are
239 two training utterances, or $R = 2$, as follows. Starting
240 from (7), we have:

$$241 \quad O(A) = \frac{w(S_1) p_A(X_1, S_1)}{\sum_{s_1} w(s_1) p_A(X_1, s_1)} + \frac{w(S_2) p_A(X_2, S_2)}{\sum_{s_2} w(s_2) p_A(X_2, s_2)}$$

$$= \frac{\sum_{s_1} w(s_1) p_A(X_1, s_1) \delta(s_1, S_1)}{\sum_{s_1} w(s_1) p_A(X_1, s_1)}$$

$$+ \frac{\sum_{s_2} w(s_2) p_A(X_2, s_2) \delta(s_2, S_2)}{\sum_{s_2} w(s_2) p_A(X_2, s_2)}$$

$$= \frac{\sum_{s_1} \sum_{s_2} w(s_1) w(s_2) p_A(X_1, s_1) p_A(X_2, s_2) [\delta(s_1, S_1) + \delta(s_2, S_2)]}{\sum_{s_1} \sum_{s_2} w(s_1) w(s_2) p_A(X_1, s_1) p_A(X_2, s_2)}$$

$$= \frac{\sum_{s_1 s_2} w(s_1 s_2) p_A(X_1, X_2, s_1, s_2) [C(s_1 s_2)]}{\sum_{s_1 s_2} w(s_1 s_2) p_A(X_1, X_2, s_1, s_2)}.$$

242 (9) 243 244 245

246 The last step used the common assumption that the train-
247 ing tokens are independent of each other. Clearly (9), is in
248 the same form of (8) when $R = 2$.

249 (2) After assuming the equivalence of (7) and (8) for
250 $R = R_0$, we now prove the equivalence for $R = R_0 +$
251 1 as follows. Again, starting from (7) for $R = R_0 +$
252 1, we have,

$$\begin{aligned}
\sum_{r=1}^{R_0+1} \frac{w(S_r)p_A(X_r, S_r)}{\sum_{s_r} w(s_r)p_A(X_r, s_r)} &= \sum_{r=1}^{R_0} \frac{w(S_r)p_A(X_r, S_r)}{\sum_{s_r} w(s_r)p_A(X_r, s_r)} + \frac{w(S_{R_0+1})p_A(X_{R_0+1}, S_{R_0+1})}{\sum_{s_{R_0+1}} w(s_{R_0+1})p_A(X_{R_0+1}, s_{R_0+1})} \\
&= \frac{\sum_{s_1 \dots s_{R_0}} w(s_1 \dots s_{R_0})p_A(X_1 \dots X_{R_0}, s_1 \dots s_{R_0})C(s_1 \dots s_{R_0})}{\sum_{s_1 \dots s_{R_0}} w(s_1 \dots s_{R_0})p_A(X_1 \dots X_{R_0}, s_1 \dots s_{R_0})} \\
&\quad + \frac{\sum_{s_{R_0+1}} w(s_{R_0+1})p_A(X_{R_0+1}, s_{R_0+1})\delta(s_{R_0+1}, S_{R_0+1})}{\sum_{s_{R_0+1}} w(s_{R_0+1})p_A(X_{R_0+1}, s_{R_0+1})} \\
&= \frac{\sum_{s_1 \dots s_{R_0}} \sum_{s_{R_0+1}} w(s_{R_0+1})p_A(X_{R_0+1}, s_{R_0+1})w(s_1 \dots s_{R_0})p_A(X_1 \dots X_{R_0}, s_1 \dots s_{R_0})[C(s_1 \dots s_{R_0}) + \delta(s_{R_0+1}, S_{R_0+1})]}{\sum_{s_1 \dots s_{R_0}} \sum_{s_{R_0+1}} w(s_{R_0+1})p_A(X_{R_0+1}, s_{R_0+1})w(s_1 \dots s_{R_0})p_A(X_1 \dots X_{R_0}, s_1 \dots s_{R_0})} \\
&= \frac{\sum_{s_1 \dots s_{R_0+1}} w(s_1 \dots s_{R_0+1})p_A(X_1 \dots X_{R_0+1}, s_1 \dots s_{R_0+1})C(s_1 \dots s_{R_0+1})}{\sum_{s_1 \dots s_{R_0+1}} w(s_1 \dots s_{R_0+1})p_A(X_1 \dots X_{R_0+1}, s_1 \dots s_{R_0+1})}, \tag{10}
\end{aligned}$$

253 that is, (8) is valid for $R = R_0 + 1$. This completes the
 254 proof by induction.

255
 256 The significance of the rational functional form of the
 257 MCE criterion is that it enables the use of the GT/EBW
 258 optimization method for discriminative training of the
 259 HMM parameters, which we elaborate below.

260 4. A new look at MCE – optimization method

261 4.1. Introduction to the growth-transformation optimization 262 technique

263 GT/EBW technique was developed for optimization of a
 264 rational function. Gopalakrishnan et al. (1991) proposed
 265 the GT/EBW based MMI estimation for the discrete
 266 HMM, and the method was extended for MMI estimation
 267 of the continuous-density HMM (CDHMM) in (Norman-
 268 din, 1991). Later Gunawardana and Byrne (2001) give an
 269 alternative method for MMI estimation of CDHMM,
 270 and its validity is proved in (Axelrod et al. (in press)). In
 271 following sections, we will present a similar method
 272 for optimization of the re-formulated MCE objective
 273 function.

274 Let $G(A)$ and $H(A)$ be two real valued functions on the
 275 parameter set A , and let the denominator function $H(A)$ be
 276 positive valued. Construct the objective function as the
 277 ratio of them to form the rational function of
 278

$$280 \quad O(A) = \frac{G(A)}{H(A)}. \tag{11}$$

281 An example of this rational function is the objective
 282 function for the MCE criterion, where
 283

$$285 \quad \begin{aligned} G(A) &= \sum_s w(s)p_A(X, s)C(s), \quad \text{and} \\ H(A) &= \sum_s w(s)p_A(X, s), \end{aligned} \tag{12}$$

and we use $s = s_1, \dots, s_R$ to denote the label sequences for
 all R training tokens, and use $X = X_1, \dots, X_R$, to denote
 the observation data sequences for all R training tokens.

As in (Gopalakrishnan et al., 1991), for the objective
 function with the form of (11), the GT-based optimization
 algorithm constructs the auxiliary function of

$$F(A; A') = G(A) - O(A')H(A) + D, \tag{13}$$

where D is a quantity independent of the parameter set A ,
 and A' denotes the parameter set obtained from the imme-
 diately previous iteration of the algorithm.

The GT algorithm starts by initializing the parameter
 set as, say, A' . (This is often accomplished by the ML
 training using, for instance, EM or Baum–Welch algo-
 rithm for HMMs.) Then, updating of the parameter set
 from A' to A proceeds by maximizing the auxiliary func-
 tion $F(A; A')$, and the process iterates until convergence is
 reached. Maximizing the auxiliary function $F(A; A')$ is
 often easier than maximizing the original rational function
 $O(A)$. It is easy to prove (Gopalakrishnan et al., 1991)
 that as long as D is a quantity not relevant to the param-
 eter set A , an increase of $F(A; A')$ guarantees an increase
 of $O(A)$.

We now define another auxiliary function from the pre-
 vious auxiliary function $F(A; A') = F(\theta)$ defined in (13).
 This new function is:

$$V(\theta; \theta') = \sum_q \int_{\chi} f(\chi, q, \theta') \log f(\chi, q, \theta) d\chi, \tag{14}$$

where the positive, real valued function $f(x, q, \theta) > 0$ is de-
 fined by

$$F(\theta) = \sum_q \int_{\chi} f(\chi, q, \theta) d\chi \tag{15}$$

and where q is a discrete variable (e.g., a state sequence in
 an HMM).

Then we have:

$$\begin{aligned}
& \log F(\theta) - \log F(\theta') \\
&= \log \frac{F(\theta)}{F(\theta')} = \log \sum_q \int_{\chi} \frac{f(\chi, q, \theta')}{F(\theta')} \frac{f(\chi, q, \theta)}{f(\chi, q, \theta')} d\chi \\
&\geq \sum_q \int_{\chi} \frac{f(\chi, q, \theta')}{F(\theta')} \log \frac{f(\chi, q, \theta)}{f(\chi, q, \theta')} d\chi \\
&= \frac{1}{F(\theta')} \left[\sum_q \int_{\chi} f(\chi, q, \theta') \log f(\chi, q, \theta) d\chi \right. \\
&\quad \left. - \sum_q \int_{\chi} f(\chi, q, \theta') \log f(\chi, q, \theta') d\chi \right] \\
&= \frac{1}{F(\theta')} [V(\theta; \theta') - V(\theta'; \theta')]. \tag{16}
\end{aligned}$$

The inequality above is due to Jensen's inequality (Jensen, 1906) applied to the concave log function. The result of (16) says that an increase in the auxiliary function $V(\theta; \theta')$ guarantees an increase in $\log F(\theta)$. Since logarithm is a monotonically increasing function, this also guarantees an increase of $F(\theta)$ and hence the original objective function $O(A)$. The technique that "transforms" the parameters from A' to A so as to increase or "grow" the values of the auxiliary functions and hence the value of the original objective function is called the growth-transformation (GT) technique (Gopalakrishnan et al., 1991). We now apply this GT technique to the Gaussian HMM with the MCE optimization criterion formulated in (8).

4.2. Application to Gaussian HMM

Substituting (12) into (13), we obtain the auxiliary function

$$\begin{aligned}
F(A; A') &= \sum_s w(s) p_A(X, s) C(s) \\
&\quad - O(A') \sum_s w(s) p_A(X, s) + D \\
&= \sum_s w(s) p_A(X, s) [C(s) - O(A')] + D \\
&= \sum_q \sum_s w(s) p_A(X, q, s) [C(s) - O(A')] + D \tag{17}
\end{aligned}$$

where q is the HMM state sequence, and $s = s_1, \dots, s_R$ is the label sequence (e.g., the word or phone sequence) for all R training tokens (including both correct or incorrect label sequences).

Follow the method used in (Gunawardana and Byrne, 2001), we can re-formulate $F(A; A')$ as follows. Let A consist of mean and variance parameters in the HMM. Since (q, s) is irrelevant with A , we have $p(X, q, s|A) = p(X|q, A)P(q, s)$, and hence

$$\begin{aligned}
F(A; A') &= \sum_q \left[\sum_s w(s) P(q, s) [C(s) - O(A')] \right] p_A(X|q) + D \\
&= \sum_q \int_{\chi} [\Gamma(A') + d(q)] p_A(\chi|q) d\chi, \tag{18}
\end{aligned}$$

where $\Gamma(A') = \delta(\chi, X) \sum_s w(s) P(q, s) [C(s) - O(A')]$, and $D = \sum_q d(q)$ is a quantity independent of the parameter set A . This quantity should be sufficiently large to guarantee that the integrand of (15) be positive, or $\Gamma(A') + d(q) > 0$ (note $p_A(\chi|q)$ in (18) is non-negative).

We now desire to construct the auxiliary function of (14) based on the auxiliary function (18). To achieve this, we first identify from (18) that

$$f(\chi, q, A) = [\Gamma(A') + d(q)] p_A(\chi|q, A), \tag{19}$$

according to (15). Then, using (14), we have

$$\begin{aligned}
V(A; A') &= \sum_q \int_{\chi} [\Gamma(A') + d(q)] p_A(\chi|q) \log \{ [\Gamma(A') + d(q)] p_A(\chi|q) \} d\chi \\
&= \sum_q \int_{\chi} [\Gamma(A') + d(q)] p_A(\chi|q) \log p_A(\chi|q) d\chi + K \\
&= \sum_q \left[\sum_s w(s) p_A(X, q, s) (C(s) - O(A')) \right] \log p_A(X|q) \\
&\quad + \sum_q d(q) \int_{\chi} p_A(\chi|q) \log p_A(\chi|q) d\chi + K. \tag{19}
\end{aligned}$$

Ignoring optimization-independent quantity K in (19), and dividing $V(A; A')$ by another optimization-independent quantity $p_A(X)$, we obtain an equivalent auxiliary function of

$$\begin{aligned}
U(A; A') &= \sum_q \left[\sum_s w(s) p_A(s|X) p_A(q|X, s) (C(s) - O(A')) \right] \log p_A(X|q) \\
&\quad + \sum_q d'(q) \int_{\chi} p_A(\chi|q) \log p_A(\chi|q) d\chi \tag{20}
\end{aligned}$$

where

$$d'(q) = d(q)/p_A(X). \tag{21}$$

Note $X = X_1, \dots, X_R$, is a large aggregate of all training data with R independent sentence tokens, and for each token $X_r = x_{r,1}, \dots, x_{r,T_r}$, the observation vector x_r , depends only on the state at time t . This enables decomposition of $\log p_A(X|q)$ and then drastic simplification of both terms in (20). To pursue the simplification, we define

$$\gamma_{i,r,s_r}(t) = p_A(q_{r,t} = i | X_r, s_r), \tag{22}$$

as the occupation probability of state i at time t , given the label sequence s_r and observation sequence X_r . Note (22) can be efficiently computed by the standard forward-backward algorithm (Rabiner and Juang, 1993). We also define

$$d(r, t, i) = \sum_{q: q_{r,t}=i} d'(q). \tag{23}$$

Then, after a series of algebraic steps, Eq. (20) can be simplified to:

$$\begin{aligned}
U(A; A') &= \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_s w(s) p_{A'}(s|X)(C(s) \\
&\quad - O(A')) \gamma_{i,r,s_r}(t) \log p_A(x_{r,t}|q_{r,t} = i) \\
&\quad + \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I d(r, t, i) \\
&\quad \times \int_{\chi_{r,t}} p_{A'}(\chi_{r,t}|q_{r,t} = i) \log p_A(\chi_{r,t}|q_{r,t} = i) d\chi_{r,t}.
\end{aligned} \tag{24}$$

402
403 We now proceed to maximize (24) with respect to the
404 Gaussian HMM's parameters, mean vectors and covari-
405 ance matrices $A = \{\mu_i, \Sigma_i\}$, $i = 1, 2, \dots, I$, in the following
406 state-conditioned Gaussian distribution:

$$408 \quad p_A(x|q = i) \propto \frac{1}{|\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right].$$

409 We set $\frac{\partial U(A; A')}{\partial A} = 0$ and solve for A given the model
410 parameters $A' = \{\mu'_i, \Sigma'_i\}$ from the previous iteration of
411 the GT/EBW. For the mean and covariance, respectively,
412 in the Gaussian at the HMM's state i , we set

$$414 \quad \frac{\partial U(A; A')}{\partial \mu_i} = 0; \quad \text{and} \quad \frac{\partial U(A; A')}{\partial \Sigma_i} = 0.$$

415 This eventually gives the ‘‘growth transformation’’ for-
416 mulas of:
417

$$419 \quad \mu_i = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t) x_t + D_i \mu'_i}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t) + D_i} \tag{25}$$

420 and
421

$$423 \quad \Sigma_i = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} [\Delta\gamma(i, r, t) (x_t - \mu_i)(x_t - \mu_i)^T] + D_i \Sigma'_i + D_i (\mu_i - \mu'_i)(\mu_i - \mu'_i)^T}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t) + D_i}, \tag{26}$$

424 where we use the new definitions of
425

$$D_i = \sum_{r=1}^R \sum_{t=1}^{T_r} d(r, t, i), \tag{27}$$

$$427 \quad \Delta\gamma(i, r, t) = \sum_s w(s) p_{A'}(s|X)(C(s) - O(A')) \gamma_{i,r,s_r}(t). \tag{28}$$

428 And we leave the detailed derivations leading to (25) and
429 (26) to the interested readers.

430 4.3. Computing $\Delta\gamma(i, r, t)$

431 The major computational steps in the above GT re-esti-
432 mation formulas are the computation of $\Delta\gamma(i, r, t)$ in (28),
433 which involves summation over all possible label sequences
434 $s = s_1, \dots, s_R$. The number of training tokens (sentence
435 strings), R , is usually very large. Hence, summation over
436 s needs to be decomposed and simplified.

437 Denote $s' = s_1, \dots, s_{r-1}$, $s'' = s_{r+1}, \dots, s_R$, $X' = X_1, \dots,$
438 X_{r-1} , $X'' = X_{r+1}, \dots, X_R$. Then, from (28), we obtain
439

$$\begin{aligned}
\Delta\gamma(i, r, t) &= \sum_{s'} \sum_{s_r} \sum_{s''} w(s', s_r, s'') p_{A'}(s', s_r, s''|X', X_r, X'') \\
&\quad \times (C(s', s_r, s'') - O(A')) \gamma_{i,r,s_r}(t) \\
&= \sum_{s_r} w(s_r) p_{A'}(s_r|X_r) \\
&\quad \times \underbrace{\left[\sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') (C(s', s_r, s'') - O(A')) \right]}_{\text{Term I}} \gamma_{i,r,s_r}(t).
\end{aligned} \tag{29} \quad 441$$

Using $C(s', s_r, s'') = C(s_r) + C(s', s'')$, Term I in (29) can
442 be simplified to
443

$$\begin{aligned}
\text{Term I} &= \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') (C(s', s_r, s'') - O(A')) \\
&= \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') C(s_r) \\
&\quad + \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') C(s', s'') \\
&\quad - O(A') \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'').
\end{aligned} \tag{30} \quad 445$$

And using
446

$$\begin{aligned}
O(A') &= \sum_{r=1}^R \frac{w(S_r) p_{A'}(X_r, S_r)}{\sum_{s_r} w(s_r) p_{A'}(X_r, s_r)} \\
&= \frac{w(S_r) p_{A'}(X_r, S_r)}{\sum_{s_r} w(s_r) p_{A'}(X_r, s_r)} + \sum_{i=1, i \neq r}^R \frac{w(S_i) p_{A'}(X_i, S_i)}{\sum_{s_i} w(s_i) p_{A'}(X_i, s_i)} \\
&= \frac{w(S_r) p_{A'}(X_r, S_r)}{\sum_{s_r} w(s_r) p_{A'}(X_r, s_r)} \\
&\quad + \frac{\sum_{s', s''} w(s', s'') p_{A'}(s', s''|X', X'') C(s', s'')}{\sum_{s', s''} w(s', s'') p_{A'}(s', s''|X', X'')} \\
&= \frac{w(S_r) p_{A'}(S_r|X_r)}{\sum_{s_r} w(s_r) p_{A'}(s_r|X_r)} \\
&\quad + \frac{\sum_{s', s''} w(s', s'') p_{A'}(s', s''|X', X'') C(s', s'')}{\sum_{s', s''} w(s', s'') p_{A'}(s', s''|X', X'')},
\end{aligned} \tag{31} \quad 448$$

we obtain:
449

$$\begin{aligned}
\text{Term I} &= \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') C(s_r) \\
&\quad + \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') C(s', s'') \\
&\quad - \frac{w(S_r) p_{A'}(S_r|X_r)}{\sum_{s_r} w(s_r) p_{A'}(s_r|X_r)} \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') \\
&\quad - \sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'') C(s', s'') \\
&= \underbrace{\sum_{s'} \sum_{s''} w(s', s'') p_{A'}(s', s''|X', X'')}_{\text{Term II}} \\
&\quad \times \left[C(s_r) - \frac{w(S_r) p_{A'}(S_r|X_r)}{\sum_{s_r} w(s_r) p_{A'}(s_r|X_r)} \right].
\end{aligned} \tag{32} \quad 451$$

The quantity above denoted by Term II can be simpli-
452 fied to:
453

$$\begin{aligned} \text{Term II} &= \frac{\sum_{s'} \sum_{s_r} \sum_{s''} w(s', s_r, s'') p_{A'}(s', s_r, s'' | X', X_r, X'')}{\sum_{s_r} w(s_r) p_{A'}(s_r | X_r)} \\ &= \frac{Q(A')}{\sum_{s_r} w(s_r) p_{A'}(s_r | X_r)} \end{aligned}$$

455

456 where we define

$$\begin{aligned} Q(A') &= \sum_s w(s) p_{A'}(s | X) \\ &= \sum_{s_1} \cdots \sum_{s_R} w(s_1) \cdots w(s_R) \cdot p_{A'}(s_1 | X_1) \cdots p_{A'}(s_R | X_R) \\ &= \prod_{r=1}^R \sum_{s_r} w(s_r) p_{A'}(s_r | X_r). \end{aligned}$$

458

459 Substituting the above terms back to, we obtain:

460

$$\begin{aligned} \Delta\gamma(i, r, t) &= \sum_{s_r} w(s_r) p_{A'}(s_r | X_r) \frac{Q(A')}{\sum_{s_r} w(s_r) p_{A'}(s_r | X_r)} \\ &\quad \times \left[C(s_r) - \frac{w(s_r) p_{A'}(s_r | X_r)}{\sum_{s_r} w(s_r) p_{A'}(s_r | X_r)} \right] \gamma_{i,r,s_r}(t). \end{aligned} \quad (30)$$

462

463 For the 1-best MCE where $w(s) \equiv 1$, $Q(A') = 1$, takes a
464 simpler form of:

$$\Delta\gamma(i, r, t) = \sum_{s_r} p_{A'}(s_r | X_r) [C(s_r) - p_{A'}(s_r | X_r)] \gamma_{i,r,s_r}(t). \quad (31)$$

466

467 4.4. Considerations for setting empirical constant D_i

468 In the GT re-estimation formulas (25) and (26), the
469 value of constant D_i is empirically set and it determines
470 the stability and convergence rate of the algorithm. We
471 now discuss the basis for setting this constant in practice.
472 From (27), (23), and (21), we have

$$\begin{aligned} D_i &= \sum_{r=1}^R \sum_{t=1}^{T_r} d(r, t, i) = \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{q, q_{r,i}=i} d'(q) \\ &= \frac{1}{p_{A'}(X)} \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{q, q_{r,i}=i} d(q). \end{aligned} \quad (32)$$

474

475 According to Jensen Inequality, the theoretical basis for
476 setting D_i is the requirement described in (18) that $d(q)$ be
477 sufficiently large to ensure that $d(q) > -\Gamma(A')$.

478 This gives

$$D_i > \frac{1}{p(X|A')} \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{q, q_{r,i}=i} -\Gamma(A'). \quad (33)$$

480

481 For the continuous density HMM case, however, $\Gamma(A')$
482 is unbounded since $\delta(\chi, X)$ is unbounded at the center point
483 $\chi = X$, and D_i needs to approach to infinite. To address this
484 issue, follow the similar derivation as in (Axelrod et al. (in
485 press)), it can be proved that with a large enough but
486 bounded D_i , the function $V(A; A')$ at (19) is still a valid
487 auxiliary function of the objective function $O(A)$, i.e.,

increasing the value of $V(A; A')$ will guarantee increase of
the value of $F(A; A')$, and so as to guarantee increase of
the value of the objective function $O(A)$.

In implementation, we are more interested in the practical
setting of D_i , which is usually determined empirically
for fast convergence. The value of D_i which we have found
practically effective for 1-best MCE is

$$D_i = E \cdot \sum_{r=1}^R p_{A'}(S_r | X_r) \sum_{s_r} p_{A'}(s_r | X_r) \sum_t \gamma_{i,r,s_r}(t), \quad (34)$$

where E is a factor controlling the learning rate. The larger
the E is, the slower the learning rate becomes, and E is usually
a factor between one and four for 1-best MCE. Extending (34) to N -best MCE, we have,

$$\begin{aligned} D_i &= E \cdot \sum_{r=1}^R \frac{Q(A')}{\sum_{s_r} w(s_r) p_{A'}(s_r | X_r)} \frac{w(S_r) p_{A'}(S_r | X_r)}{\sum_{s_r} w(s_r) p_{A'}(s_r | X_r)} \\ &\quad \times \sum_{s_r} w(s_r) p_{A'}(s_r | X_r) \sum_t \gamma_{i,r,s_r}(t). \end{aligned} \quad (35)$$

5. Use of Lattice for representing competitive candidates in MCE training

In the above novel development of the MCE training
technique, N -best lists are used to represent the competing
candidates for discriminative learning. In many speech recognition
tasks, in order to make the competing candidates sufficiently rich,
 N in the N -best lists needs to be very large (e.g., in the order
of millions). This will create computational difficulties. To overcome
such difficulties, we can use a lattice to serve as a compressed form
of a very large N -best list, which has been shown to be successful
and critical in MMI and MPE learning (Woodland and Povey,
2000; Povey, 2004). The previous work that discussed the use of
lattices for MCE training was reported in (Schluter et al., 2001),
where the misclassification measure takes the approximate form of

$$d_r(X_r, A) = -\log p_A^z(X_r, S_r) + \log \sum_{s_r, s_r \neq S_r} p_A^z(X_r, s_r), \quad (36)$$

to the misclassification measure of (2). This is a special case
of our approximate form of (3) where $w(s_r) \equiv 1$ for all strings
including the incorrect (competing) strings $\{s_r: s_r \neq S_r\}$, and
the correct string S_r .

In the above special case and with $\alpha = 1$, our earlier
results in (30) and (35) become simplified to

$$\Delta\gamma(i, r, t) = \sum_{s_r} p_{A'}(s_r | X_r) [C(s_r) - p_{A'}(s_r | X_r)] \gamma_{i,r,s_r}(t), \quad (37)$$

$$D_i = E \cdot \sum_{r=1}^R p_{A'}(S_r | X_r) \sum_t \sum_{s_r} p_{A'}(s_r | X_r) \gamma_{i,r,s_r}(t). \quad (38)$$

In this section, instead of computing $\Delta\gamma(i, r, t)$ of (37)
and in (38) a brute-force manner by summing a huge number
of strings of s_r for the very large N -best list as represented
by a lattice, we use an approximation that makes

536 the computation of (37) and (38) practically feasible. This
 537 then gives a solution for lattice-based MCE parameter esti-
 538 mation after using the computed results of (37) and (38) in
 539 the estimation formulas (25) and (26). This solution differs
 540 markedly from that reported in (Schluter et al., 2001). Spe-
 541 cifically, our solution does not require removing the correct
 542 word string S_r from the lattice. In contrast, removal of S_r in
 543 the lattice is required by the solution provided in (Schluter
 544 et al., 2001), which is more difficult to implement in prac-
 545 tice than our solution. In addition, our solution has the the-
 546 oretical appeal of guaranteed algorithm convergence since
 547 it is derived based on GT/EBW for a rational function.
 548 The solution provided in (Schluter et al., 2001) does not
 549 have such convergence guarantee.

550 To compute (37), we first use $C(s_r) = \delta(s_r, S_r)$ for MCE
 551 to rewrite (37) as
 552

$$\begin{aligned} \Delta\gamma(i, r, t) &= \sum_{s_r} p_{A'}(s_r | X_r) C(s_r) \gamma_{i,r,s_r}(t) \\ &\quad - \sum_{s_r} p_{A'}(s_r | X_r) p_{A'}(S_r | X_r) \gamma_{i,r,s_r}(t) \\ &= p_{A'}(S_r | X_r) \gamma_{i,r,S_r}(t) - p_{A'}(S_r | X_r) \\ &\quad \times \underbrace{\sum_{s_r} p_{A'}(s_r | X_r) \gamma_{i,r,s_r}(t)}_{\mathcal{Y}}. \end{aligned} \quad (39)$$

555 In (39), since the correct string S_r is known, $\gamma_{i,r,S_r}(t)$ can
 556 be computed straightforwardly by the standard forward-
 557 backward algorithm for the HMM (Rabiner and Juang,
 558 1993). The main computation thus lies in
 559

$$T = \sum_{s_r} p_{A'}(s_r | X_r) \gamma_{i,r,s_r}(t) \quad (40)$$

562 and

$$p_{A'}(S_r | X_r) = \frac{p_{A'}(S_r, X_r)}{p_{A'}(X_r)}. \quad (41)$$

566 for $\Delta\gamma(i, r, t)$ of (37), as well as for D_i in (38).

567 To efficiently compute (40) and (41) for the lattice repre-
 568 sentation of strings of s_r , we need to make a mild approx-
 569 imation. A lattice is a compact representation of a large list
 570 of strings. It is an acyclic directed graph consisting of a
 571 number of nodes and a set of directed arcs each connecting
 572 two nodes. A typical arc is denoted as q , and an arc corre-
 573 sponds to a substring (e.g., a word in a sentence). Two time
 574 stamps, b_q and e_q , are associated with each arc, providing
 575 an estimate of the segment boundaries for the substring.
 576 For a time slice t within the arc segment q , we have
 577 $b_q \leq t \leq e_q$.

578 We will show below that (40) and (41) can both be com-
 579 puted efficiently by a forward-backward algorithm on the
 580 lattice after the mild assumption that HMM state
 581 sequences are independent across arcs, that is,

$$\gamma_{i,r,q}(t) \approx \gamma_{i,r,s_r}(t) \quad \text{when } b_q \leq t \leq e_q \text{ and } q \in s_r. \quad (42)$$

This approximation says that within the segment of arc
 q , its occupancy posterior probability $\gamma_{i,r,s_r}(t) = p(q_{r,t} =$
 $i | X_r, s_r, A')$ given the observation sequence X_r for the sen-
 tence-level string s_r that passes arc q approximates the pos-
 terior probability $\gamma_{i,r,q}(t)$ for arc q alone. The justification of
 approximation (42) is that the state sequence within arc
 q should be roughly independent of other arcs. This was
 called ‘‘exact-matching’’ approximation in (Povey, 2004).
 To see this, let s_r be composed of three sub-strings: s'_r , q ,
 s''_r . Then the right hand side of (42) can be analyzed to be

$$\begin{aligned} \gamma_{i,r,s_r}(t : b_q \leq t \leq e_q) &= p_{A'}(q_{r,t:b_q \leq t \leq e_q} = i | X_r, s_r) \\ &= p_{A'}(q_{r,t:b_q \leq t \leq e_q} = i | X_r, s'_r, q, s''_r) \\ &\approx p_{A'}(q_{r,t:b_q \leq t \leq e_q} = i | X_r, q) \end{aligned} \quad 596$$

597 which is the left hand side of (42).

598 The essence of approximation (42) is to decouple the
 599 dependency on the local arc q from the entire string s_r . This
 600 enables drastic simplification of the computation in (40)
 601 and (41), which we discuss below.

602 As we discussed earlier, the principal computation bur-
 603 den in (40) is the huge number (N) of summation terms for
 604 s_r for the equivalent N -best list of a lattice. Using approx-
 605 imation of (42), we can drastically reduce the computation
 606 by the following simplification:
 607

$$\begin{aligned} T &\approx \sum_{s_r} p_{A'}(s_r | X_r) \gamma_{i,r,q}(t) \\ &= \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot \sum_{s_r: q \in s_r} p_{A'}(s_r | X_r) \\ &= \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot p_{A'}(q | X_r) \\ &= \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot \frac{p_{A'}(q, X_r)}{p_{A'}(X_r)}. \end{aligned} \quad (43) \quad 609$$

610 The key quantities in (43) can be efficiently computed as
 611 follows (See the derivation in Appendix I):
 612

$$p_{A'}(q, X_r) = \alpha(q) \beta(q); \quad (44)$$

$$p_{A'}(X_r) = \sum_{q:q \in \{\text{ending arcs}\}} p_{A'}(q, X_r) = \sum_{q:q \in \{\text{ending arcs}\}} \alpha(q) \quad (45) \quad 614$$

615 where the ‘‘forward’’ and ‘‘backward’’ probabilities are de-
 616 fined by

$$\alpha(q) = \sum_{p(\text{preceding } q)} p_{A'}(p, q, X'_r(q), X_r(q)) = p_{A'}(q, X'_r(q), X_r(q)); \quad (46)$$

$$\beta(q) = \sum_{v(\text{succeeding } q)} p_{A'}(v, X''_r(q) | q) = p_{A'}(X''_r(q) | q). \quad (47) \quad 619$$

620 In (46), $X'_r(q)$ denotes the r th training token’s observa-
 621 tion sequence preceding arc q , i.e., during $1 \leq t < b_q$.
 622 $X_r(q)$ is the observation sequence bounded by arc q with
 623 $b_q \leq t \leq e_q$. $X''_r(q)$ in (47) denotes the observation sequence
 624 succeeding arc q , or during $e_q < t \leq T_r$.

625 For each arc q in the lattice, $\alpha(q)$ and $\beta(q)$ can be
 626 computed by the following efficient forward and back-

ward recursions, respectively (See the derivation in Appendix II):

$$\alpha(q) = \sum_{p(\text{preceding } q)} p_{A'}(q|p)p_{A'}(X_r(q)|q)\alpha(p) \quad (48)$$

and

$$\beta(q) = \sum_{v(\text{succeeding } q)} p_{A'}(v|q)p_{A'}(X_r(v)|v)\beta(v), \quad (49)$$

where $\alpha(q)$ is initialized at the starting arc q_0 by $\alpha(q_0) = \pi(q_0)p_{A'}(X_r(q_0)|q_0)$, and $\beta(q)$ initialized at the ending arc q_E by $\beta(q_E) = 1$.

Using forward probability $\alpha(q)$, we can efficiently compute (41) as follows. Since $p(X_r|A') = \sum_{q,q \in \{\text{ending arcs}\}} \alpha(q)$ and $p_{A'}(S_r, X_r) = p_{A'}(X_r|S_r)p(S_r)$, we have

$$p(S_r|X_r, A') = \frac{p(X_r|S_r, A')p(S_r)}{\sum_{q,q \in \{\text{ending arcs}\}} \alpha(q)}. \quad (50)$$

6. Summary and discussion

HMMs are continuing to play a central role in speech recognition research and technology deployment, where training techniques for the HMM parameters have been a critical determinant for the speech recognition accuracy and user satisfaction level. While discriminative training for HMMs, typified by the MCE technique, has been pursued with a relatively long history, it is not until recently that the traditional gradient-based MCE optimization technique has been questioned (Macherey et al., 2005; He and Chou, 2003). In this paper, we provide a fresh look at the MCE technique not only from the perspective of the optimization technique, but also of the objective function. The key technical contribution of this paper is the establishment of a non-trivial framework in which the MCE objective function is re-formulated as a rational function for multiple sentence-level training tokens. And we show that the N -best representation of the competitive candidates in MCE training amounts to a special weighting function in the newly formulated MCE objective function. As a consequence of this re-formulation, we most naturally derive the new optimization method for discriminatively estimating HMM parameters based on GT/EBW. This method has been successfully implemented in a speech recognition system, and the positive experimental results can be found in (He et al., 2006).

In addition to the usual treatment of MCE training using the N -best paradigm, in this paper, we also provide further, more difficult technical detail for the use of lattices as a richer representation of competing candidates. This treatment can be considered as a technical guide for implementing MCE training in large-scale speech recognition systems. We are currently experimenting with this approach.

Appendix I. Derivation of Eqs. (44) and (45)

Given the “forward” and “backward” probabilities defined as (46) and (47), as well as $X'_r(q)$, $X_r(q)$ and $X''_r(q)$ defined in Section 5, a derivation of (44) and (45) is provided below.

Derivation of (44):

$$\begin{aligned} p_{A'}(q, X_r) &= \sum_{p(\text{preceding } q)} \sum_{v(\text{succeeding } q)} p_{A'}(p, q, v, X'_r(q), X_r(q), X''_r(q)) \\ &= \sum_{p(\text{preceding } q)} \sum_{v(\text{succeeding } q)} p_{A'}(p, q, X'_r(q), X_r(q)) p_{A'} \\ &\quad \times (v, X''_r(q)|p, q, X'_r(q), X_r(q)) \\ &= \sum_{p(\text{preceding } q)} p_{A'}(p, q, X'_r(q), X_r(q)) \\ &\quad \times \sum_{v(\text{succeeding } q)} p_{A'}(v, X''_r(q)|q) = \alpha(q)\beta(q). \end{aligned}$$

Derivation of (45):

$$\begin{aligned} p_{A'}(X_r) &= \sum_{q,q \in \{\text{ending arcs}\}} p_{A'}(q, X_r) \\ &= \sum_{q,q \in \{\text{ending arcs}\}} p_{A'}(q, X'_r(q), X_r(q)) \\ &= \sum_{q,q \in \{\text{ending arcs}\}} \alpha(q). \end{aligned}$$

Appendix II. Derivation of (48) and (49) for the forward–backward computation

Given the “forward” and “backward” probabilities defined as (46) and (47), as well as $X'_r(q)$, $X_r(q)$ and $X''_r(q)$ defined in Section 5, a derivation of (48) and (49) is provided here.

Forward computation for $\alpha(q)$:

$$\begin{aligned} \alpha(q) &= \sum_{p(\text{preceding } q)} p_{A'}(p, q, X'_r(q), X_r(q)) \\ &= \sum_{p(\text{preceding } q)} p_{A'}(p, q, X'_r(p), X_r(p), X_r(q)) \\ &= \sum_{p(\text{preceding } q)} p_{A'}(q, X_r(q)|p, X'_r(p), X_r(p)) p_{A'}(p, X'_r(p), X_r(p)) \\ &= \sum_{p(\text{preceding } q)} p_{A'}(q|p)p_{A'}(X_r(q)|q)p_{A'}(p, X'_r(p), X_r(p)) \\ &= \sum_{p(\text{preceding } q)} p_{A'}(q|p)p_{A'}(X_r(q)|q)\alpha(p). \end{aligned}$$

Backward computation for $\beta(q)$:

$$\begin{aligned} \beta(q) &= \sum_{v(\text{succeeding } q)} p_{A'}(v, X''_r(q)|q) \\ &= \sum_{v(\text{succeeding } q)} p_{A'}(v, X_r(v), X''_r(v)|q) \\ &= \sum_{v(\text{succeeding } q)} p_{A'}(v, X_r(v)|q)p_{A'}(X''_r(v)|q, v, X_r(v)) \end{aligned}$$

$$\begin{aligned}
 &= \sum_{v(\text{succeeding } q)} p_{A'}(v|q)p_{A'}(X_r(v)|q, v)p_{A'}(X_r''(v)|v) \\
 &= \sum_{v(\text{succeeding } q)} p_{A'}(v|q)p_{A'}(X_r(v)|v)\beta(v)
 \end{aligned}$$

703

704 **References**

705 Axelrod, S., Goel, V., Gopinath, R., Olsen, P., Visweswariah, K., in press.
 706 Discriminative Estimation of Subspace Constrained Gaussian Mixture
 707 Models for Speech Recognition, IEEE Trans. Audio Speech Language
 708 Process., <http://ieeexplore.ieee.org/iel5/10376/32978/101109TASL>
 709 [2006872617.pdf](http://ieeexplore.ieee.org/iel5/10376/32978/101109TASL).
 710 Bahl, L., Jelinek, F., Mercer, R., 1987. A Maximum likelihood approach
 711 to continuous speech recognition. IEEE Trans. Pattern Anal. Mach.
 712 Intell. PAMI-5, 179–190.
 713 Brown, P., 1987. The Acoustic Modeling Problem in Automatic Speech
 714 Recognition, Ph.D. thesis, Carnegie Mellon University.
 715 Chou, W., 2003. Minimum classification error approach in pattern
 716 recognition. In: Chou, W., Juang, B.-H. (Eds.), Pattern Recognition in
 717 Speech and Language Processing. CRC Press, pp. 1–49.
 718 Deng, L., Wu, J., Droppo, J., Acero, A., 2005. Analysis and comparison
 719 of two feature extraction/compensation algorithms. IEEE Signal
 720 Process. Lett. 12 (6), 477–480.
 721 Deng, L., Yu, D., Acero, A., 2005. A generative modeling framework for
 722 structured hidden speech dynamics. In: Proc. of Neural Information
 723 Processing System (NIPS) Workshop, Whistler, BC, Canada, Decem-
 724 ber 2005.
 725 Deng, L., O'Shaughnessy, D., 2003. SPEECH PROCESSING – A
 726 Dynamic and Optimization-Oriented Approach. Marcel Dekker Inc.,
 727 New York, NY, USA.
 728 Gopalakrishnan, P., Kanevsky, D., Nadas, A., Nahamoo, D., 1991. An
 729 inequality for rational functions with applications to some statistical
 730 estimation problems. IEEE Trans. Inf. Theory. 37, 107–113.
 731 Gunawardana, A., Byrne, W., 2001. Discriminative speaker adaptation
 732 with conditional maximum likelihood linear regression. In: Proc.
 733 EUROSPEECH.
 734 He, X., Chou, W., 2003. Minimum classification error linear regression for
 735 acoustic model adaptation of continuous density HMMs. In: Proc.
 736 ICASSP.
 737 He, X., Deng, L., Chou, W., 2006. A novel learning method for hidden
 738 Markov models in speech and audio processing. In: Proc. IEEE
 739 Workshop on Multimedia Signal Processing, Victoria, BC.

Jensen, J.L.W.V., 1906. Sur les fonctions convexes et les inegalites entre les
 740 valeurs moyennes. Acta Math., 175–193. 741
 Juang, B.-H., Katagiri, S., 1992. Discriminative learning for minimum
 742 error classification. IEEE Trans. Signal Process. 40 (12), 3043–3054. 743
 Juang, B.-H., Chou, W., Lee, C.-H., 1997. Minimum classification error
 744 rate methods for speech recognition. IEEE Trans. Speech Audio
 745 Process. 5. 746
 Macherey, W., Haferkamp, L., Schluter, R., Ney, H., 2005. Investigations
 747 on error minimizing training criteria for discriminative training in
 748 automatic speech Recognition. In: Proc. Interspeech, Lisbon, Portugal,
 749 pp. 2133–2136. 750
 McDermott, E., Hazen, T., Roux, J., Nakamura, A., Katagiri, S., in press.
 751 Discriminative training for large vocabulary speech recognition using
 752 minimum classification error. IEEE Trans. Audio Speech Language
 753 Process., <http://www.kecl.ntt.co.jp/icl/signal/erik/index-j.htm>. 754
 Normandin, Y., 1991. Hidden Markov Models, Maximum Mutual
 755 Information Estimation, and the Speech Recognition Problem, Ph.D.
 756 dissertation, McGill University, Montreal. 757
 Povey, D., 2004. Discriminative Training for Large Vocabulary
 758 Speech Recognition, Ph.D. thesis, Cambridge University, Cam-
 759 bridge, UK. 760
 Povey, D., Gales, M.J.F., Kim, D.Y., Woodland, P.C., 2003. MMI-MAP
 761 and MPE-MAP for acoustic model adaptation, In: Proc. Eurospeech.
 762
 Povey, D., Kingsbury, B., Mangu, L., Saon, G., Soltau, H., Zweig, G.,
 763
 2004. fMPE: Discriminatively trained features for speech recognition.
 764
 In: Proc. DARPA EARS RT-04 Workshop, November 7–10, Pali-
 765 sades, NY, Paper No. 35. 766
 Povey, D., Woodland, P.C., 2002. Minimum phone error and I-Smooth-
 767 ing for improved discriminative training. In: Proc. ICASSP. 768
 Rabiner, L., Juang, B.-H., 1993. Fundamentals of Speech Recognition.
 769 Prentice Hall, Englewood Cliffs, New Jersey. 770
 Rathinavalu, C., Deng, L., 1998. Speech trajectory discrimination using
 771 the minimum classification error learning. IEEE Trans. Speech Audio
 772 Process. 6 (6), 505–515. 773
 Roux, J., McDermott, E., 2005. Optimization for discriminative training.
 774
 In: Proc. INTERSPEECH. 775
 Schluter, R., Macherey, W., Muller, B., Ney, H., 2001. Comparison of
 776 discriminative training criteria and optimization methods for speech
 777 recognition. Speech Commun. 34, 287–310. 778
 Woodland, P.C., Povey, D., 2000. Large scale discriminative training for
 779 speech recognition. In: Proc. ITRW ASR, ISCA. 780
 781