

# A STUDY OF N-GRAM AND DECISION TREE LETTER LANGUAGE MODELING METHODS

Gerasimos Potamianos<sup>a,†</sup>, Frederick Jelinek<sup>b</sup>

<sup>a</sup>Speech and Image Processing Services Research Lab, AT&T Labs–Research, Florham Park, NJ 07932-0971, USA

<sup>b</sup>Center for Language and Speech Processing, The Johns Hopkins University, Baltimore, MD 21218, USA

Received 11 August 1997; revised 2 March 1998

## Abstract

The goal of this paper is to investigate various language model smoothing techniques and decision tree based language model design algorithms. For this purpose, we build language models for printable characters (letters), based on the Brown corpus. We consider two classes of models for the text generation process: The  $n$ -gram language model and various decision tree based language models. In the first part of the paper, we compare the most popular smoothing algorithms applied to the former. We conclude that the bottom-up deleted interpolation algorithm performs the best in the task of  $n$ -gram letter language model smoothing, significantly outperforming the back-off smoothing technique for large values of  $n$ . In the second part of the paper, we consider various decision tree development algorithms. Among them, a  $K$ -means clustering type algorithm for the design of the decision tree questions gives the best results. However, the  $n$ -gram language model outperforms the decision tree language models for letter language modeling. We believe that this is due to the predictive nature of letter strings, which seems to be naturally modeled by  $n$ -grams.

## Zusammenfassung

Das Ziel dieses Beitrags ist verschiedene Techniken zur Glättung von Sprachmodellen und Algorithmen zum Entwurf von Sprachmodellen auf der Basis von Entscheidungsbäumen zu untersuchen. Zu diesem Zweck verwenden wir den Brown-Korpus um Modelle von Buchstabenfolgen zu erstellen. Wir betrachten zwei Klassen von Modellen zur Textgenerierung: das  $n$ -Gramm Sprachmodell sowie verschiedene auf Entscheidungsbäumen basierende Verfahren. Im ersten Teil dieses Beitrags vergleichen wir die am häufigsten benutzten Glättungsalgorithmen angewandt auf  $n$ -Gramme. Wir folgern, daß der "bottom-up deleted interpolation"-Algorithmus am besten zur Glättung von  $n$ -Gramm Sprachmodellen geeignet ist und für große  $n$  dem "back-off"-Verfahren deutlich überlegen ist. Im zweiten Teil dieses Beitrags betrachten wir dann verschiedene Algorithmen zur Bildung von Entscheidungsbäumen. Unter diesen erzielt ein  $K$ -means-ähnlicher Algorithmus die besten Ergebnisse beim Entwurf der Fragen die der Entscheidungsbaum stellt. Für die Modellierung

von Buchstabenfolgen erzielt das  $n$ -Gramm Sprachmodell aber trotzdem noch bessere Ergebnisse als alle Entscheidungsbäume. Wir glauben, daß dies durch die Fähigkeit der  $n$ -Gramme Buchstabenfolgen vorherzusagen begründet ist.

## Résumé

Le but de cet article est d'étudier différentes techniques de lissage de modèles de langage et différents algorithmes de construction de modèles de langage à base d'arbres de décision. Pour cela, nous construisons des modèles de langage pour des caractères écrits (lettres) à partir du Brown corpus. Nous considérons deux classes de modèles pour le processus de génération du texte : le modèle de langage  $n$ -gram, et différents modèles de langage à base d'arbres de décision. Dans la première partie de l'article, nous comparons les algorithmes de lissage les plus couramment appliqués au modèle de langage  $n$ -gram. L'algorithme "bottom-up deleted interpolation" donne les meilleurs résultats pour le lissage du modèle de langage  $n$ -gram, dépassant de façon significative la technique de lissage "back-off" pour de grandes valeurs de  $n$ . Dans la seconde partie de l'article, nous considérons différents algorithmes de développement d'arbres de décision. Parmi eux, un algorithme de type classification  $K$ -means aboutit aux meilleurs résultats pour la construction des arbres de décision. Cependant, le modèle de langage  $n$ -gram fournit de meilleurs résultats que les modèles de langage à arbres de décision pour la modélisation du langage écrit. Nous croyons que cela est dû à la nature prédictive des chaînes de caractères, qui semble être naturellement modélisée par les  $n$ -grams.

*Keywords:* Language modeling;  $n$ -grams; Decision trees; Smoothing; Laws of succession; Back-off language model; Deleted interpolation; Brown corpus

## 1. Introduction

Although less interesting than *word language models*, *letter language models* are suitable for certain applications, such as text compression (Bell et al., 1990), spelling letter automatic recognition (Betz and Hild, 1995), and estimation of language entropy (Brown et al., 1992a; Shannon,

<sup>†</sup> Corresponding author. E-mail: makis@research.att.com.

1951). In addition, letter language models are simpler to construct, since their *vocabulary* size is significantly smaller than that of a typical word language model. In fact, their small vocabulary size allows a clear demonstration of the effects of *sparseness* of the training data on the language model performance. In an  $n$ -gram letter language model, data sparseness appears gradually, as  $n$  increases. In contrast, it abruptly arises in large vocabulary word language models, unless more sophisticated *class* language models are used (Brown et al., 1992b). Clearly, letter language models provide simple testing grounds of the language modeling techniques that constitute the subject of this paper. Such techniques are purely statistical methods, therefore, their generalization to word language modeling is straightforward.

The most widely adopted language model is the  $n$ -gram language model (Jelinek, 1997). Our motivation has been to outperform it by exploring *decision tree* (Breiman et al., 1984) based language models. Our results however are negative: Our decision tree based letter language models fail to improve the  $n$ -gram one, unless certain  $n$ -gram like constraints are imposed on the decision tree question design. Nevertheless, this paper contains a number of original results: First, although decision tree based language models have appeared before (Bahl et al., 1989), our work constitutes the first fair comparison between them and  $n$ -gram language models. In addition, the employment of the *K-means clustering* algorithm for decision tree question design (Chou, 1991) is its first application to language modeling. Finally, our decision tree *smoothing* algorithm is worth reporting: We propose the use of either *bottom-up deleted interpolation* (Jelinek and Mercer, 1980) or *back-off* smoothing (Katz, 1987), depending on the *depth* of each tree *leaf*. The investigation of suitable decision tree smoothing techniques has led us to a comparative study of the most popular smoothing techniques. Our conclusion is important: The bottom-up deleted interpolation smoothing is the most suitable smoothing technique for letter language models, being significantly more robust to sparseness of the training data than both *top-down* deleted interpolation and back-off smoothing.

Let  $\omega_1^m = \omega_1\omega_2\omega_3\dots\omega_m$  be a sequence of  $m$  characters of text, where  $\omega_i \in \mathcal{V}$ , for  $i = 1, 2, 3, \dots, m$ , and  $\mathcal{V}$  denotes the vocabulary, which contains  $|\mathcal{V}|$  distinct characters. We consider the text generation process as a random one, therefore,  $\omega_1^m$  is drawn from an unknown probability mass function  $\overline{\text{Pr}}[\omega_1^m] = \prod_{i=1}^m \overline{\text{Pr}}[\omega_i | \omega_1^{i-1}]$ . We are interested in letter language models that approximate these unknown conditional probability mass functions by *homogeneous* (i.e., independent of the character text location  $i$ ) conditional probabilities

$$\begin{aligned} \overline{\text{Pr}}[\omega_i | \omega_1^{i-1}] &\simeq \text{Pr}[\omega_i | \Phi(\omega_{i-n+1}^{i-1})] \\ &= \text{Pr}[\omega_0 | \Phi(\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1})] > 0, \end{aligned} \quad (1)$$

for all  $\omega_{i-n+1}^i \in \mathcal{V}^n$ , and for a given positive integer  $n$ . In (1),  $\Phi$  denotes an *equivalence classification* on  $\mathcal{V}^{n-1}$ , whereas,  $\omega_{-j} \in \mathcal{V}$  denotes the character lying  $j$  positions to the left of character  $\omega_0 \in \mathcal{V}$ , that it helps to predict. Clearly, two issues are of interest in the design of language model (1): First, the specification of the equivalence classification  $\Phi$ , and, second, the determination of probabilities

(1). In this work, we consider  $n$ -gram and decision tree based language models, as a means of specifying  $\Phi$ , and various smoothing techniques, as a means of determining the conditional probability mass functions (1).

In this case study for letter language modeling, all algorithms discussed are compared on the *Brown corpus* (Kucera and Francis, 1967), a popular corpus with the language modeling community (Brown et al., 1992a; Chen and Goodman, 1996; Ristad and Thomas, 1995, 1997a, 1997b). More specifically, the Brown corpus, consisting of approximately 6.1 million characters with a vocabulary of 79 characters (see Table 2), is randomly divided into three sets, namely a *development*, a *held-out* (or, *check*), and a *test* set of approximately 4.9, 0.7, and 0.5 million characters, respectively<sup>1</sup> (the union of the development and held-out sets constitutes the training set). Counts  $C(\omega_0, \omega_{-1}, \dots, \omega_{-n+1})$ , for all  $n$ -tuples  $(\omega_0, \omega_{-1}, \dots, \omega_{-n+1}) \in \mathcal{V}^n$ , are collected from the three sets, denoted by  $C_d(\bullet)$ ,  $C_c(\bullet)$ , and  $C_t(\bullet)$ , respectively. Given a language model (1), we define the *minus log-probability per character* on each of the three sets, denoted by  $\mathcal{LP}_d$ ,  $\mathcal{LP}_c$ , and  $\mathcal{LP}_t$ , respectively, as

$$\begin{aligned} \mathcal{LP}_s = & \quad (2) \\ & \frac{\sum_{(\omega_0, \omega_{-1}, \dots, \omega_{-n+1}) \in \mathcal{V}^n} C_s(\omega_0, \omega_{-1}, \dots, \omega_{-n+1}) \log_2 \text{Pr}[\omega_0 | \Phi(\omega_{-1}, \dots, \omega_{-n+1})]}{\sum_{(\omega_0, \omega_{-1}, \dots, \omega_{-n+1}) \in \mathcal{V}^n} C_s(\omega_0, \omega_{-1}, \dots, \omega_{-n+1})}, \end{aligned}$$

for  $s = d, c, t$ , and measured in bits per character. Quantity (2) is consistently used throughout this paper, as a means of language model design and assessment. In more detail, counts  $C_d(\bullet)$  are used to obtain the equivalence classification  $\Phi$ , as well as, the *maximum likelihood* estimates of probabilities (1), by minimizing  $\mathcal{LP}_d$ . In addition, counts  $C_c(\bullet)$  are used to optimize language model parameters, related to the smoothing of these maximum likelihood estimates, in order to improve the language model generalization to the held-out set. Minimization of  $\mathcal{LP}_c$  is employed for this task. Finally, the performance of the resulting language models is assessed by means of  $\mathcal{LP}_t$ . With no loss of generality, alternatives to the held-out data scheme, described above, such as the *bootstrap*, the *leave-one-out*, and the *N-fold cross validation* methods (Efron, 1982; Ney et al., 1995) are not considered here. Such methods exhibit well known statistical advantages over the held-out data approach, however, they are significantly more computationally expensive.

This paper is structured as follows: In Section 2, various  $n$ -gram language model smoothing algorithms are discussed, namely various *laws of succession*<sup>2</sup> (Bell et al., 1990; Ney et al., 1995; Ristad, 1995), the *back-off* method (Katz, 1987; Ney et al., 1995), and *deleted interpolation* (Bahl et al., 1983, 1991; Chen and Goodman, 1996; Jelinek and Mercer, 1980; Ristad and Thomas, 1997b). Results on the

<sup>1</sup>All held-out and test set characters appear in the development set, therefore, we have not introduced an *out of vocabulary* character in  $\mathcal{V}$ .

<sup>2</sup>We use this terminology to refer to smoothing techniques that take into account development set event occurrences only within the equivalence class  $\Phi(\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1})$  under consideration.

Brown corpus are reported in Section 2.4. Section 3 is devoted to decision tree based language models. Two decision tree language model development algorithms are presented in Sections 3.2 and 3.3. Decision tree smoothing is discussed in Section 3.4, whereas, results on the Brown corpus are reported in Section 3.5. Section 4 summarizes our conclusions. The reader is encouraged to consult our technical report for a more detailed version of this work (Potamianos and Jelinek, 1997).<sup>3</sup>

## 2. A comparison of various language model smoothing techniques for the $n$ -gram language model

A conceptually simple but surprisingly powerful language model is the  $n$ -gram language model that corresponds to the equivalence classification (see also (1))  $\Phi(\mathbf{h}_{n-1}) = \mathbf{h}_{n-1}$ , where  $\mathbf{h}_{n-1} = (\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1})$  denotes the *history* of depth  $n-1$ , preceding character  $\omega_0$ . We are interested in estimating  $\Pr[\omega_0 = v | \mathbf{h}_{n-1}]$ , for all  $v \in \mathcal{V}$ , and all histories  $\mathbf{h}_{n-1} \in \mathcal{V}^{n-1}$ . In practice, at least some such histories never occur in the development set. In these cases, we set

$$\Pr[v | \omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1}] = \Pr[v | \omega_{-1}, \omega_{-2}, \dots, \omega_{-K}],$$

for all  $v \in \mathcal{V}$ , where

$$K = \max_{1 \leq k < n} \{k : C_d(\omega_{-1}, \omega_{-2}, \dots, \omega_{-k}) > 0\}.$$

Therefore, in the  $n$ -gram language model, we need estimates of  $\Pr[v | \mathbf{h}_k]$ , whenever  $C_d(\mathbf{h}_k) > 0$ , and for all<sup>4</sup>  $k = 0, 1, 2, \dots, n-1$ . By minimizing  $\mathcal{L}\mathcal{P}_d$ , we obtain maximum likelihood estimates of such probabilities, given by (Jelinek, 1997)

$$\Pr_{\text{ML}}[v | \mathbf{h}_k] = \frac{C_d(v, \mathbf{h}_k)}{C_d(\mathbf{h}_k)}, \text{ for all } v \in \mathcal{V}. \quad (3)$$

The resulting (based on (3)) minus log-probability on the development set,  $\mathcal{L}\mathcal{P}_d$ , equals its *entropy* (Shannon, 1951), which we denote by  $\mathcal{H}_d$ . Notice that, in contrast to the decision tree based language models, discussed in Section 3, minimization of  $\mathcal{L}\mathcal{P}_d$  (or, equivalently, of  $\mathcal{H}_d$ , given (3)) is not explicitly used to obtain the  $n$ -gram language model equivalence classification  $\Phi$ . Of course,  $\mathcal{H}_d$  is a decreasing function of  $n$  (see also Table 1).

Clearly, smoothing of the probability mass functions (3) is necessary, in order to provide the desired language model generalization to unseen data. Three classes of such smoothing algorithms are discussed next. Experimental results on the Brown corpus are presented in Section 2.4.

### 2.1. Laws of succession

Laws of succession constitute a group of smoothing techniques, where probability mass functions  $\Pr[v | \mathbf{h}_k]$ , with  $0 \leq k < n$  and  $C_d(\mathbf{h}_k) > 0$ , are estimated only by means of counts  $C_d(v, \mathbf{h}_k)$ , for all  $v \in \mathcal{V}$ , as well as, by means of certain parameters that are either a-priori chosen, or optimized by taking into account held-out set event counts. Therefore, no information is used from a coarser partitioning of the development set history space.

A variety of laws of succession exist in the literature (Bell et al., 1990; Good, 1953; Ney et al., 1995; Ristad, 1995). We have carried out a detailed comparison of a number of them, in the case study of letter  $n$ -gram language modeling on the Brown corpus (Potamianos and Jelinek, 1997). The laws of succession listed below consistently outperformed the rest. We denote by  $\mathbf{h}$  any history  $\mathbf{h}_k$  with  $C_d(\mathbf{h}) > 0$ , where  $0 \leq k < n$ , and by  $q(\mathbf{h}) = |\{v \in \mathcal{V} : C_d(v, \mathbf{h}) > 0\}|$ .

- Sharpened uniform cardinality prior based law (Ristad, 1995), given by

$$\Pr_{\text{LOS.1}}[v | \mathbf{h}] = \begin{cases} \frac{C_d(v, \mathbf{h})}{C_d(\mathbf{h})}, & \text{if } q(\mathbf{h}) = |\mathcal{V}|; \text{ otherwise:} \\ \frac{C_d(v, \mathbf{h})}{C_d(\mathbf{h})} \frac{C_d(\mathbf{h}) [C_d(\mathbf{h}) + 1] + q(\mathbf{h}) [1 - q(\mathbf{h})]}{[C_d^2(\mathbf{h}) + C_d(\mathbf{h}) + 2q(\mathbf{h})]}, & \text{if } C_d(v, \mathbf{h}) > 0, \\ \frac{q(\mathbf{h}) [q(\mathbf{h}) + 1]}{[|\mathcal{V}| - q(\mathbf{h})] [C_d^2(\mathbf{h}) + C_d(\mathbf{h}) + 2q(\mathbf{h})]}, & \text{if } C_d(v, \mathbf{h}) = 0. \end{cases} \quad (4.1)$$

- Absolute discounting by half law (Bell et al., 1990; Ney et al., 1995), given by

$$\Pr_{\text{LOS.2}}[v | \mathbf{h}] = \begin{cases} \frac{C_d(v, \mathbf{h})}{C_d(\mathbf{h})}, & \text{if } q(\mathbf{h}) = |\mathcal{V}|; \text{ otherwise:} \\ \frac{C_d(v, \mathbf{h}) - 0.5}{C_d(\mathbf{h})}, & \text{if } C_d(v, \mathbf{h}) > 0, \\ \frac{0.5 q(\mathbf{h})}{[|\mathcal{V}| - q(\mathbf{h})] C_d(\mathbf{h})}, & \text{if } C_d(v, \mathbf{h}) = 0. \end{cases} \quad (4.2)$$

- Absolute discounting law (Ney et al., 1995), which is a generalization of (4.2) using a discounting constant  $0 < \delta(\mathbf{h}) < 1$ ,

$$\Pr_{\text{LOS.3}}[v | \mathbf{h}] = \begin{cases} \frac{C_d(v, \mathbf{h})}{C_d(\mathbf{h})}, & \text{if } q(\mathbf{h}) = |\mathcal{V}|; \text{ otherwise:} \\ \frac{C_d(v, \mathbf{h}) - \delta(\mathbf{h})}{C_d(\mathbf{h})}, & \text{if } C_d(v, \mathbf{h}) > 0, \\ \frac{\delta(\mathbf{h}) q(\mathbf{h})}{[|\mathcal{V}| - q(\mathbf{h})] C_d(\mathbf{h})}, & \text{if } C_d(v, \mathbf{h}) = 0. \end{cases} \quad (4.3)$$

In (4.3), the smoothing coefficients  $\delta(\mathbf{h})$  have to be obtained for the histories  $\mathbf{h}$  of interest. It is reasonable to choose these coefficients to minimize the minus log-probability on the held-out set,  $\mathcal{L}\mathcal{P}_e$  (see (2)), hoping that this will provide good generalization of the resulting language model

<sup>3</sup>This report can be downloaded from <http://www.research.att.com/~makis/CLSP97.13.ps>.

<sup>4</sup>The case  $k = 0$  corresponds to the uni-gram. In such a case,  $C_d(v, \mathbf{h}_0) = C_d(v)$ .

to the unseen test set (i.e., result in a small  $\mathcal{LP}_t$ ). However, we need to guard against over-fitting of these parameters to the held-out data set. We can meet both goals by *bucketing* the histories  $\mathbf{h}$  of interest in terms of a chosen feature, for example, the development set history count  $C_d(\mathbf{h})$ , as suggested by Bahl et al. (1983), while simultaneously assuring that enough held-out set events are contained within the constructed history *buckets*. All histories within a bucket are assigned a single value of the smoothing constant. Features alternative to  $C_d(\mathbf{h})$  can also be used to construct history buckets. Indeed, the use of  $C_d(\mathbf{h})/q(\mathbf{h})$  is advocated by Chen and Goodman (1996), whereas, the two-dimensional feature  $(C_d(\mathbf{h}), q(\mathbf{h}))$  is used by Ristad and Thomas (1997b). In the following, we briefly discuss our bucketing algorithm in terms of feature  $C_d(\mathbf{h})$  (Potamianos and Jelinek, 1997), and its application to estimating  $\delta(\mathbf{h})$  in (4.3). Extensions to other history bucketing features can be made in a straightforward manner.

Let  $\mathcal{C}$  denote the set of histories  $\mathbf{h}$  of interest. For example (see (4.3)),  $\mathcal{C} = \{\mathbf{h}_k: C_d(\mathbf{h}_k) > 0, \text{ and } q(\mathbf{h}_k) < |\mathcal{V}|\}$ . We now seek integer constants  $B_0, B_1, \dots, B_L$ , that partition the set  $\mathcal{C}$  into  $L$  sets (buckets)

$$S_j = \{\mathbf{h} \in \mathcal{C}: B_{j-1} \leq C_d(\mathbf{h}) < B_j\}, \text{ for } j = 1, 2, \dots, L. \quad (5)$$

Such constants are chosen to satisfy

$$\begin{aligned} \min_{\mathbf{h} \in \mathcal{C}} \{C_d(\mathbf{h})\} &= B_0 < B_1 < B_2 < \dots \\ &< B_{L-1} < B_L = \max_{\mathbf{h} \in \mathcal{C}} \{C_d(\mathbf{h})\} + 1, \end{aligned} \quad (6a)$$

$$\sum_{\mathbf{h} \in S_j} C_c(\mathbf{h}) \geq C_{\min}, \quad (6b)$$

and, approximately (subject to (6a) and (6b)),<sup>5</sup>

$$B_j \simeq \lceil b B_{j-1} \rceil, \quad (6c)$$

for  $j = 1, 2, \dots, L$ . In (6),  $b > 1$  is a real constant that controls the “width” of buckets (5), whereas,  $C_{\min}$  denotes the minimum required number of held-out set events in bucket  $S_j$  to ensure reliable estimation of tied  $\delta_j = \delta(\mathbf{h})$ , for  $\mathbf{h} \in S_j$ ,  $j = 1, 2, \dots, L$ .<sup>6</sup> The tied smoothing coefficients are then determined by minimizing  $\mathcal{LP}_c$  (equivalently, by maximizing  $-\mathcal{LP}_c$ ), as (see (2) and (4.3))

$$\begin{aligned} \hat{\delta}_j &= \arg \max_{0 < \delta_j < 1} \left\{ \sum_{\mathbf{h} \in S_j} \sum_{v \in \mathcal{V}: C_d(v, \mathbf{h}) > 0} C_c(v, \mathbf{h}) \log_2 \frac{C_d(v, \mathbf{h}) - \delta_j}{C_d(\mathbf{h})} \right. \\ &\quad \left. + \sum_{\mathbf{h} \in S_j} \sum_{v \in \mathcal{V}: C_d(v, \mathbf{h}) = 0} C_c(v, \mathbf{h}) \log_2 \frac{\delta_j q(\mathbf{h})}{[|\mathcal{V}| - q(\mathbf{h})] C_d(\mathbf{h})} \right\}, \end{aligned} \quad (7a)$$

for  $j = 1, 2, \dots, L$ . Optimization (7a) can be simply carried out by means of the *Newton-Raphson* algorithm (Press

et al., 1988), or the simpler, but slower, *bisection* method (Press et al., 1988), for solving

$$\begin{aligned} \sum_{\mathbf{h} \in S_j} \sum_{v \in \mathcal{V}: C_d(v, \mathbf{h}) > 0} C_c(v, \mathbf{h}) [\delta_j - C_d(v, \mathbf{h})]^{-1} \\ + \sum_{\mathbf{h} \in S_j} \sum_{v \in \mathcal{V}: C_d(v, \mathbf{h}) = 0} C_c(v, \mathbf{h}) \delta_j^{-1} = 0, \end{aligned} \quad (7b)$$

in the compact set  $[\epsilon, 1 - \epsilon]$ , where  $\epsilon$  is a small positive number, such as  $\epsilon = 10^{-5}$ . Notice that there exists a unique solution to (7b), because its left hand side is the first derivative of the argument of (7a), which is a concave function with respect to  $\delta_j$ .

## 2.2. The back-off approach

The main mechanism of the back-off smoothing algorithm (Katz, 1987; Ney et al., 1995) is the discounting of the maximum likelihood estimates (3) of seen event probabilities  $\text{Pr}_{\text{ML}}[v | \mathbf{h}_k]$  by means of any law of succession and the assignment of the remaining probability mass to the unseen events (such that,  $C_d(v, \mathbf{h}_k) = 0$  and  $C_d(\mathbf{h}_k) > 0$ ), not in a uniform way (as is the case with any law of succession), but according to the back-off smoothed probability mass function of the immediate ancestor history  $\mathbf{h}_{k-1}$  of the history  $\mathbf{h}_k$  under consideration.

The back-off smoothed estimates of  $\text{Pr}[v | \mathbf{h}_k]$ , for all  $v \in \mathcal{V}$ , whenever  $C_d(\mathbf{h}_k) > 0$ , with  $0 \leq k \leq n-1$ , are defined by the initialization<sup>7</sup>

$$\text{Pr}_{\text{BOF}}[v | \mathbf{h}_0] = \text{Pr}_{\text{LOS}}[v | \mathbf{h}_0], \quad (8a)$$

and the top-down recursion, for  $i = 1, 2, \dots, k$ ,

$$\begin{aligned} \text{Pr}_{\text{BOF}}[v | \mathbf{h}_i] \\ = \begin{cases} \text{Pr}_{\text{LOS}}[v | \mathbf{h}_i], & \text{if } q(\mathbf{h}_i) = |\mathcal{V}|; \text{ otherwise:} \\ \text{Pr}_{\text{LOS}}[v | \mathbf{h}_i], & \text{if } C_d(v, \mathbf{h}_i) > 0, \\ \beta(\mathbf{h}_i) \text{Pr}_{\text{BOF}}[v | \mathbf{h}_{i-1}], & \text{if } C_d(v, \mathbf{h}_i) = 0, \end{cases} \end{aligned} \quad (8b)$$

where, if  $q(\mathbf{h}_i) < |\mathcal{V}|$ ,

$$\beta(\mathbf{h}_i) = \frac{1 - \sum_{v \in \mathcal{V}: C_d(v, \mathbf{h}_i) > 0} \text{Pr}_{\text{LOS}}[v | \mathbf{h}_i]}{\sum_{v \in \mathcal{V}: C_d(v, \mathbf{h}_i) = 0} \text{Pr}_{\text{BOF}}[v | \mathbf{h}_{i-1}]}, \quad (8c)$$

which guarantees that (8b) is a valid probability mass function. In (8), any law of succession (4) could be used. We denote the resulting back-off probability mass function accordingly; for example, by  $\text{Pr}_{\text{BOF},1}[v | \mathbf{h}_k]$ , if  $\text{Pr}_{\text{LOS},1}[v | \mathbf{h}_k]$  is employed. Clearly, when using back-off smoothing in combination with (4.3), we need to estimate the smoothing parameters  $\delta(\mathbf{h}_i)$ . It is not hard to see that (7b) is applicable here as well. However, (6) and (7b) should be consecutively applied to seen histories of depths  $0, 1, \dots, n-1$ , in a top-down fashion (Potamianos and Jelinek, 1997).

<sup>7</sup>In (8), we denote by  $\mathbf{h}_{i-1}$  the immediate ancestor of history  $\mathbf{h}_i$ ; for example, if  $\mathbf{h}_i = (\omega_{-1}, \omega_{-2}, \dots, \omega_{-i})$ , then  $\mathbf{h}_{i-1} = (\omega_{-1}, \omega_{-2}, \dots, \omega_{-i+1})$ . Notice that  $\mathbf{h}_0$  corresponds to the unigram.

<sup>5</sup>We denote  $\lceil x \rceil = \min\{i: i \text{ is integer, and } i \geq x\}$ .

<sup>6</sup>Throughout our language modeling experiments on the Brown corpus, and whenever smoothing algorithms require bucketing for optimally determining smoothing coefficients, we have used  $C_{\min} = \min\{\lceil |\mathcal{V}|/4 \rceil, \sum_{\mathbf{h} \in \mathcal{C}} C_c(\mathbf{h})\}$  and  $b = 1.2$ .

### 2.3. Deleted interpolation

Let us again consider any history  $\mathbf{h}_k$ , such that  $C_d(\mathbf{h}_k) > 0$  and  $0 \leq k \leq n-1$ . We are interested in determining  $\Pr[v | \mathbf{h}_k]$ , for all  $v \in \mathcal{V}$ . For such a history, we consider the series of its ancestors  $\mathbf{h}_i$ , and maximum likelihood estimates  $\Pr_{\text{ML}}[v | \mathbf{h}_i]$  (see (3)), for  $i = k-1, k-2, \dots, 1, 0$ . In addition, we consider the “null” history  $\mathbf{h}_{-1} = \emptyset$ , with  $\Pr_{\text{ML}}[v | \mathbf{h}_{-1}] = 1 / |\mathcal{V}|$ . The deleted interpolation algorithm (Bahl et al., 1983; Jelinek and Mercer, 1980), estimates the desired  $\Pr[v | \mathbf{h}_k]$ ,  $v \in \mathcal{V}$ , as

$$\Pr_{\text{DI}}[v | \mathbf{h}_k] = \sum_{i=-1}^k \lambda^{(i)}(\mathbf{h}_k) \Pr_{\text{ML}}[v | \mathbf{h}_i], \quad (9a)$$

where  $\lambda^{(i)}(\mathbf{h}_k)$  are smoothing coefficients that satisfy

$$0 < \lambda^{(i)}(\mathbf{h}_k) < 1, \quad \text{for } i = -1, 0, 1, 2, \dots, k, \\ \text{and } \sum_{i=-1}^k \lambda^{(i)}(\mathbf{h}_k) = 1. \quad (9b)$$

The smoothing coefficients should be optimally chosen to minimize the minus log-probability on the held-out data set,  $\mathcal{L}\mathcal{P}_c$  (see (2)). In this paper, we consider two algorithms for estimating these coefficients, namely the *top-down* and the *bottom-up* deleted interpolation algorithms. We denote the resulting probability mass functions by  $\Pr_{\text{DI.TD}}[v | \mathbf{h}_k]$  and  $\Pr_{\text{DI.BU}}[v | \mathbf{h}_k]$ , respectively.

#### 2.3.1. Top-down deleted interpolation

The top-down deleted interpolation algorithm sequentially derives  $\Pr_{\text{DI.TD}}[v | \mathbf{h}_0]$ ,  $\Pr_{\text{DI.TD}}[v | \mathbf{h}_1]$ ,  $\dots$ ,  $\Pr_{\text{DI.TD}}[v | \mathbf{h}_k]$ , for all histories  $\mathbf{h}_k$  of interest, as illustrated in Fig. 1a. In detail, let us define sets

$$\mathcal{C}_i = \{\mathbf{h}_i \in \mathcal{V}^i : C_d(\mathbf{h}_i) > 0\}, \quad \text{for } i = 0, 1, \dots, n-1, \quad (10a)$$

(clearly,  $\mathcal{C}_0$  is a singleton), and apply (6) to obtain history buckets

$$S_{i,j} = \{\mathbf{h}_i \in \mathcal{C}_i : B_{i,j-1} \leq C_d(\mathbf{h}_i) < B_{i,j}\}, \\ \text{for } j = 1, 2, \dots, L_i, \quad i = 0, 1, \dots, n-1. \quad (10b)$$

Then, the top-down deleted interpolation algorithm consists of the initialization

$$\Pr_{\text{DI.TD}}[v | \mathbf{h}_{-1}] = \Pr_{\text{ML}}[v | \mathbf{h}_{-1}] = \frac{1}{|\mathcal{V}|}, \quad (11a)$$

and the following two steps, iterated over all buckets  $S_{i,j}$ ,  $j = 1, 2, \dots, L_i$ , and over all levels  $i = 0, 1, \dots, n-1$ : The optimization step (see also (2))

$$\hat{\lambda}_{i,j} = \arg \max_{0 < \lambda < 1} \left\{ \sum_{\mathbf{h}_i \in S_{i,j}} \sum_{v \in \mathcal{V}} C_c(v, \mathbf{h}_i) \times \log_2 \left[ \lambda \Pr_{\text{DI.TD}}[v | \mathbf{h}_{i-1}] + (1-\lambda) \Pr_{\text{ML}}[v | \mathbf{h}_i] \right] \right\}, \quad (11b)$$

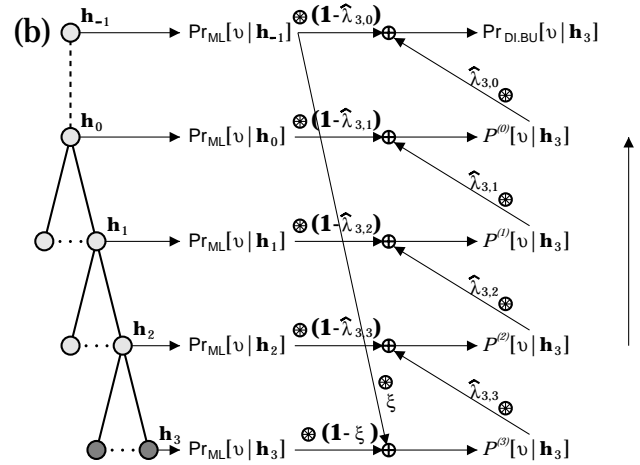
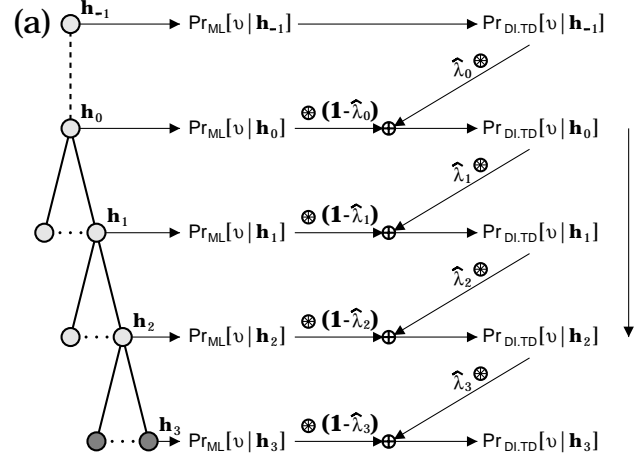


Figure 1: Demonstration of the two deleted interpolation algorithms, discussed in Section 2.3, in the case  $k = 3$ . (a): Top-down algorithm. (b): Bottom-up algorithm.

and the substitution step

$$\Pr_{\text{DI.TD}}[v | \mathbf{h}_i] \\ = \hat{\lambda}_{i,j} \Pr_{\text{DI.TD}}[v | \mathbf{h}_{i-1}] + (1 - \hat{\lambda}_{i,j}) \Pr_{\text{ML}}[v | \mathbf{h}_i], \quad (11c)$$

where, in (11c),  $\mathbf{h}_i \in S_{i,j}$ . Optimization (11b) can be easily carried out, by means of the Newton-Raphson algorithm for solving

$$\frac{\sum_{\mathbf{h}_i \in S_{i,j}} \sum_{v \in \mathcal{V}} C_c(v, \mathbf{h}_i) \times \Pr_{\text{DI.TD}}[v | \mathbf{h}_{i-1}] - \Pr_{\text{ML}}[v | \mathbf{h}_i]}{\lambda \left[ \Pr_{\text{DI.TD}}[v | \mathbf{h}_{i-1}] - \Pr_{\text{ML}}[v | \mathbf{h}_i] \right] + \Pr_{\text{ML}}[v | \mathbf{h}_i]} = 0, \quad (11d)$$

in the compact set  $[\epsilon, 1-\epsilon]$  (e.g.,  $\epsilon = 10^{-5}$ ). Alternatively, the bisection method can be used to solve (11d) (Bahl et al., 1991). Similarly to solving (7b), solving (11d) is tractable, because its left hand side is the first derivative of the argument of (11b), which is a concave function with respect to

$\lambda$ . The final smoothing coefficients of (9) are given by

$$\lambda^{(i)}(\mathbf{h}_k) = \begin{cases} \prod_{k'=0}^k \widehat{\lambda}_{k', j_{k'}}, & \text{if } i = -1, \\ (1 - \widehat{\lambda}_{i, j_i}) \prod_{k'=i+1}^k \widehat{\lambda}_{k', j_{k'}}, & \text{if } 0 \leq i \leq k, \end{cases} \quad (11e)$$

where  $j_{k'} = j$ , whenever  $\mathbf{h}_{k'} \in S_{k', j}$ , where  $0 \leq k' \leq k$  (Potamianos and Jelinek, 1997).

### 2.3.2. Bottom-up deleted interpolation

The bottom-up deleted interpolation algorithm constructs a series of estimates of probability mass function  $\text{Pr}_{\text{DI.BU}}[v | \mathbf{h}_k]$ ,  $v \in \mathcal{V}$ , for all histories  $\mathbf{h}_k$  that satisfy  $C_d(\mathbf{h}_k) > 0$  and  $0 \leq k \leq n-1$ . Let us denote such estimates by  $P^{(i)}[v | \mathbf{h}_k]$ ,  $v \in \mathcal{V}$ , for  $i = k, k-1, \dots, 0, -1$ . The initial estimate  $P^{(k)}[v | \mathbf{h}_k]$  of  $\text{Pr}_{\text{DI.BU}}[v | \mathbf{h}_k]$  is a slightly smoothed version of the maximum likelihood estimate  $\text{Pr}_{\text{ML}}[v | \mathbf{h}_k]$  at level  $k$  (see (12a), below), whereas, the final estimate  $P^{(-1)}[v | \mathbf{h}_k]$  is the desired smoothed probability mass function  $\text{Pr}_{\text{DI.BU}}[v | \mathbf{h}_k]$  (see also Fig. 1b).

In detail, suppose that, for some  $0 \leq k \leq n-1$ , we wish to estimate  $\text{Pr}_{\text{DI.BU}}[v | \mathbf{h}_k]$ , whenever  $C_d(\mathbf{h}_k) > 0$ . First, we bucket the histories in set  $C_k$  (see (10a)) into sets (buckets)  $S_{k, j}$ ,  $j = 1, 2, \dots, L_k$  (see (10b)), by using (6). Then, the bottom-up deleted interpolation algorithm consists of the initialization step

$$P^{(k)}[v | \mathbf{h}_k] = (1 - \xi) \text{Pr}_{\text{ML}}[v | \mathbf{h}_k] + \xi \text{Pr}_{\text{ML}}[v | \mathbf{h}_{-1}], \quad (12a)$$

where  $0 < \xi < 1$  is a sufficiently small positive parameter, chosen in an ad-hoc manner (e.g.,  $\xi = 10^{-5}$ ), and the following two steps, iterated over all buckets  $S_{k, j}$ ,  $j = 1, 2, \dots, L_k$ , and over all  $i = k, k-1, \dots, 1, 0$ : The optimization step (see also (2))

$$\widehat{\lambda}_{k, i, j} = \arg \max_{0 < \lambda < 1} \left\{ \sum_{\mathbf{h}_k \in S_{k, j}} \sum_{v \in \mathcal{V}} C_c(v, \mathbf{h}_k) \times \log_2 \left[ \lambda P^{(i)}[v | \mathbf{h}_k] + (1 - \lambda) \text{Pr}_{\text{ML}}[v | \mathbf{h}_{i-1}] \right] \right\}, \quad (12b)$$

performed by means of the Newton-Raphson algorithm, or the bisection method (similarly to the solution of (7b) and (11d)), as a solution of

$$\frac{\sum_{\mathbf{h}_k \in S_{k, j}} \sum_{v \in \mathcal{V}} C_c(v, \mathbf{h}_k) \times P^{(i)}[v | \mathbf{h}_k] - \text{Pr}_{\text{ML}}[v | \mathbf{h}_{i-1}]}{\lambda \left[ P^{(i)}[v | \mathbf{h}_k] - \text{Pr}_{\text{ML}}[v | \mathbf{h}_{i-1}] \right] + \text{Pr}_{\text{ML}}[v | \mathbf{h}_{i-1}]} = 0,$$

in the compact set  $[\epsilon, 1 - \epsilon]$  (e.g.,  $\epsilon = 10^{-5}$ ), and the substitution step

$$P^{(i-1)}[v | \mathbf{h}_k] = \widehat{\lambda}_{k, i, j} P^{(i)}[v | \mathbf{h}_k] + (1 - \widehat{\lambda}_{k, i, j}) \text{Pr}_{\text{ML}}[v | \mathbf{h}_{i-1}], \quad (12c)$$

where, in (12c),  $\mathbf{h}_k \in S_{k, j}$ . Finally,

$$\text{Pr}_{\text{DI.BU}}[v | \mathbf{h}_k] = P^{(-1)}[v | \mathbf{h}_k], \quad \text{for all } v \in \mathcal{V}. \quad (12d)$$

The final smoothing coefficients of (9) are given by

$$\lambda^{(i)}(\mathbf{h}_k) = \begin{cases} 1 - \widehat{\lambda}_{k, 0, j_k} + \prod_{k'=0}^{k+1} \widehat{\lambda}_{k', k', j_k}, & \text{if } i = -1, \\ (1 - \widehat{\lambda}_{k, i+1, j_k}) \prod_{k'=0}^i \widehat{\lambda}_{k', k', j_k}, & \text{if } 0 \leq i \leq k, \end{cases} \quad (12e)$$

where  $j_k = j$ , whenever  $\mathbf{h}_k \in S_{k, j}$ , and  $\widehat{\lambda}_{k, k+1, j_k} = \xi$ . Notice that (12) have to be calculated for all  $k = 0, 1, \dots, n-1$  (Potamianos and Jelinek, 1997).

Obviously, the top-down deleted interpolation algorithm is computationally less intensive than the bottom-up one. This is due to the fact that the former naturally determines  $\text{Pr}_{\text{DI.TD}}[v | \mathbf{h}_{i-1}]$ , before determining  $\text{Pr}_{\text{DI.TD}}[v | \mathbf{h}_i]$ . However, the latter algorithm gives rise to superior language models, possibly because the history bucketing, as well as, successive optimizations (12b) always occur at the finest level of interest  $k$ .

Notice that both top-down and bottom-up deleted interpolation algorithms avoid the simultaneous optimization of all smoothing coefficients, and, therefore, they are sub-optimal. Alternatively, one can choose to simultaneously estimate all  $k+2$  such coefficients in (9a), by maximizing the held-out data log-probability,  $-\mathcal{L}\mathcal{P}_c$ , over the simplex (9b). This is a concave function with respect to the vector of smoothing coefficients (Bahl et al., 1991), therefore, a unique solution to the problem can be found by means of a *conjugate gradient* method (Press et al., 1988). The desired smoothing coefficients can be tied as in the bottom-up deleted interpolation algorithm. However, our preliminary experiments indicate that this approach results in insignificant improvement (over the use of (12)) in both held-out and test set log-probabilities, and at a significantly higher computational cost, as  $k$  becomes larger. Alternatively, one can use the *expectation-maximization* algorithm to simultaneously estimate the smoothing coefficients in (9) (Jelinek and Mercer, 1980; Ristad and Thomas, 1997b). However, this method has been reported to be significantly more computationally expensive than optimization (11b), even in the simple case  $k = 0$ , considered by Bahl et al. (1991). In addition, Ristad and Thomas (1997b) compare the use of the expectation-maximization algorithm to our bottom-up deleted interpolation technique for  $n$ -gram language model smoothing, and they report no significant performance difference between the two methods.

## 2.4. Experimental results

We now compare the performance of the three laws of succession (4.1)-(4.3), their three back-off variants (8), and the two deleted interpolation smoothing algorithms (11), (12), in the case of the  $n$ -gram letter language model on the Brown corpus, for  $n = 1, 2, \dots, 10$ . In Table 1, we depict the test set minus log-probability,  $\mathcal{L}\mathcal{P}_t$ , of the resulting language models, as well as, the development set entropy,  $\mathcal{H}_d$ .

$n$	$\mathcal{H}_d$	LOS.1	LOS.2	LOS.3	BOF.1	BOF.2	BOF.3	DI.TD	DI.BU
1	4.378	4.380	4.380	4.380	4.380	4.380	4.380	4.380	4.380
2	3.469	3.471	3.471	3.470	3.470	3.471	3.470	3.470	3.470
3	2.847	2.851	2.852	2.851	2.851	2.851	2.850	2.851	2.850
4	2.288	2.333	2.334	2.331	2.327	2.327	2.324	2.328	2.326
5	1.883	2.038	2.040	2.034	2.007	2.008	2.002	2.016	2.007
6	1.609	<b>1.979</b>	<b>1.983</b>	<b>1.975</b>	1.884	1.886	1.876	1.894	1.878
7	1.375	2.071	2.074	2.065	<b>1.864</b>	<b>1.867</b>	<b>1.854</b>	1.853	1.831
8	1.135	2.233	2.236	2.227	1.888	1.891	1.877	1.837	1.811
9	0.899	2.401	2.403	2.394	1.927	1.930	1.915	1.828	1.801
10	<b>0.687</b>	2.554	2.556	2.539	1.971	1.974	1.948	<b>1.824</b>	<b>1.796</b>

Table 1: Minus log-probability on our Brown corpus test subset,  $\mathcal{L}\mathcal{P}_t$ , of the  $n$ -gram language model, smoothed by the algorithms discussed in Section 2, for  $n = 1, 2, \dots, 10$ . The development set entropy,  $\mathcal{H}_d$ , is also listed. All results are measured in bits per character

As expected, the laws of succession perform the worst: The resulting  $n$ -gram models are over-trained for  $n > 6$ . The absolute discounting law (LOS.3) is slightly better than the other laws of succession, since it allows estimation of the smoothing coefficients  $\delta(\mathbf{h})$  on a held-out set. The back-off smoothing technique performs significantly better than the laws of succession. Back-off smoothing based on the absolute discounting law (BOF.3) is slightly better than the rest (BOF.1 and BOF.2). However, the resulting  $n$ -gram models are over-trained for  $n > 7$ . On the other hand, both deleted interpolation algorithms exhibit robustness to sparseness of the training data, with the bottom-up algorithm achieving smaller minus log-probability on the test set.

Our best  $n$ -gram language model achieves a minus log-probability of 1.796 bits per character on the Brown corpus test set, and it is obtained by smoothing the 10-gram language model by means of the bottom-up deleted interpolation algorithm. This result is better than the 1.91 bits per character achieved by the context model of Ristad and Thomas (1995), but worse than the 1.75 bits per character reported by Brown et al. (1992a). There, a sophisticated tri-gram word language model with dedicated spelling, case, and spacing sub-models is employed, and, moreover, 600 times more training data, obtained from various corpora, is used. Finally, the hierarchical non-emitting Markov language model, introduced by Ristad and Thomas (1997a, 1997b), achieves 1.73 bits per character on the same task. This model does not satisfy (1), and it is shown to be strictly more powerful than the  $n$ -gram language model. We believe that part of its superior performance is due to the use of  $N$ -fold cross validation, re-estimation of (3) on basis of the entire training set, and employment of a two-dimensional feature history bucketing scheme, when estimating necessary smoothing coefficients (Ristad and Thomas, 1997b).

## 2.5. Remarks on smoothing algorithms

In Section 2, we presented a variety of smoothing algorithms for language models. Among them, we concluded that the bottom-up deleted interpolation is the most suitable for our  $n$ -gram letter language models on the Brown corpus, being very robust to data sparseness. Both back-off smoothing algorithms and laws of succession fail to exhibit

this robustness, with the back-off smoothing algorithm being significantly better than its corresponding law of succession, especially for sparse data. Interestingly, the back-off scheme based on absolute discounting compares favorably to the bottom-up deleted interpolation algorithm for values of  $n \leq 6$ , thus being a viable alternative to it. However, when the training data is really sparse ( $n \geq 7$ ), one has to resort to the safety of the bottom-up deleted interpolation algorithm, which significantly outperforms the back-off smoothing algorithm.

These conclusions are in line with experimental results reported in Chen and Goodman (1996). There, and for three word language modeling tasks, the back-off smoothed word  $n$ -gram language models are consistently outperformed in sparse training data domains by the deleted interpolation smoothed ones (word tri-gram language models, trained on “insufficient” data). This is typically reversed in the case of bi-gram word language models, where, clearly, the training data domain is not as sparse. A possible cause for the inferior performance of back-off smoothing in sparse data domains is the fact that (4.3) and (8b) do not discount the development set counts  $C_d(v, \mathbf{h})$  enough (recall that  $0 < \delta(\mathbf{h}) < 1$ ), thus, favoring higher order history events over lower order ones. A remedy to this problem could be to discount  $C_d(v, \mathbf{h})$  by  $\delta(v, \mathbf{h})$ , instead of  $\delta(\mathbf{h})$ , in (4.3), where  $0 < \delta(v, \mathbf{h}) < C_d(v, \mathbf{h})$ . Such discounting schemes have not been considered in this paper. It is worth mentioning that our history bucketing scheme, as implemented in this work, has improved the performance of back-off smoothing based on the absolute discounting law, over the implementation reported by Ney et al. (1995). There, for each history level  $k$ , all discounting coefficients are tied to a single value; i.e.,  $\delta_k = \delta(\mathbf{h}_k)$ , for all  $\mathbf{h}_k$  with  $C_d(\mathbf{h}_k) > 0$  and  $0 \leq k \leq n-1$  (see also (4.3)). Similarly to our experimental results, Ney et al. (1995) report that the absolute discounting law is the most appropriate for back-off smoothing.

So far, our exposition has been in reference to the  $n$ -gram language model. However, the above smoothing algorithms are readily applicable to more general language models. Indeed, in Section 3.4, we discuss bottom-up deleted interpolation and back-off smoothing for decision tree language models.

### 3. Decision tree based language modeling

Decision trees constitute popular classifiers (Breiman et al., 1984). Since the language modeling problem (1) is a classification problem, it is natural to investigate decision tree based language models. Such models have been developed before (Bahl et al., 1989), however, our work is novel in the following aspects: (a): We apply  $K$ -means clustering for decision tree *question* design (Chou, 1991); (b): We compare this algorithm to one based on a *binary encoding* of the vocabulary (Jelinek, 1997); (c): We propose a robust decision tree smoothing algorithm; and (d): We provide a fair comparison between decision tree and  $n$ -gram language models. Experimental results on the Brown corpus are presented in Section 3.5.

#### 3.1. Decision tree notation and design issues

A tree  $\mathbf{T}$  consists of a finite collection of *nodes*  $\mathbf{T} = \{t_0, t_1, \dots\}$ , where  $t_0$  denotes the tree *root*. Every node has a unique *parent* node, denoted by  $PAR(t) \in \mathbf{T}$  (except for  $PAR(t_0) = NULL$ ), and  $K(t) \geq 0$  *children* nodes, denoted by  $KID_i(t) \in \mathbf{T}$ , for  $i = 1, 2, \dots, K(t)$ . For every node  $t \in \mathbf{T}$ , we recursively define its *depth*, as

$$DEPTH(t) = \begin{cases} 0, & \text{if } t = t_0; \\ \text{otherwise:} \\ DEPTH(PAR(t)) + 1, & \end{cases} \quad (13a)$$

and the set of its *ancestors*  $ANC(t) = \{ANC_i(t) : i = -1, 0, 1, \dots, DEPTH(t)\}$ , as

$$ANC_i(t) = \begin{cases} t, & \text{if } i = DEPTH(t); \\ \text{otherwise:} \\ PAR(ANC_{i+1}(t)), & \\ \text{for } i = DEPTH(t) - 1, \dots, 1, 0, -1. & \end{cases} \quad (13b)$$

Clearly,  $ANC_0(t) = t_0$ , and  $ANC_{-1}(t) = NULL$ . We denote the set of tree *leaves* by  $\mathbf{L} = \{t \in \mathbf{T} : K(t) = 0\} \subseteq \mathbf{T}$ . Here, we exclusively consider *binary* and *ternary* trees that satisfy  $K(t) = 2$  and  $K(t) = 3$ , for all  $t \in \mathbf{T} - \mathbf{L}$ , respectively. In addition, our ternary trees satisfy  $KID_3(t) \in \mathbf{L}$ , for all  $t \in \mathbf{T} - \mathbf{L}$ . Nodes  $KID_3(t)$  are called *middle nodes*, and their set is denoted by  $\mathbf{M} \subset \mathbf{L} \subseteq \mathbf{T}$  (see also Fig. 2).

A decision tree is a tree  $\mathbf{T}$  with a question  $Q(t)$  assigned at each node  $t \in \mathbf{T} - \mathbf{L}$ , and a probability mass function  $\Pr[v|t]$ , assigned at each leaf  $t \in \mathbf{L}$ . Question  $Q(t)$  partitions node  $t$  into its  $K(t)$  children nodes, according to a mapping (question)  $Q(t) : t \rightarrow \{KID_1(t), \dots, KID_{K(t)}(t)\}$ . In the context of language modeling,  $v \in \mathcal{V}$ , every node  $t \in \mathbf{T}$  corresponds to a cluster of histories  $\Phi_t$  (see also (1)), whereas, every question  $Q(t)$  partitions  $\Phi_t$  into  $K(t)$  new history clusters  $\Phi_{KID_i(t)}$ , for  $i = 1, 2, \dots, K(t)$ . Given the tree  $\mathbf{T}$ , we define (compare to (3))

$$f(v|t) = \Pr_{ML}[v|t] = \frac{C_d(v, t)}{C_d(t)},$$

and  $f(t|t') = \frac{C_d(t)}{C_d(t')},$  (14)

for all  $v \in \mathcal{V}$ ,  $t \in \mathbf{T}$ ,  $t' \in ANC(t)$ , and  $t, t' \neq NULL$ . In (14),  $C_d(v, t) = \sum_{\mathbf{h}_{n-1} \in \Phi_t} C_d(v, \mathbf{h}_{n-1})$ , and it denotes the total number of development set occurrences of

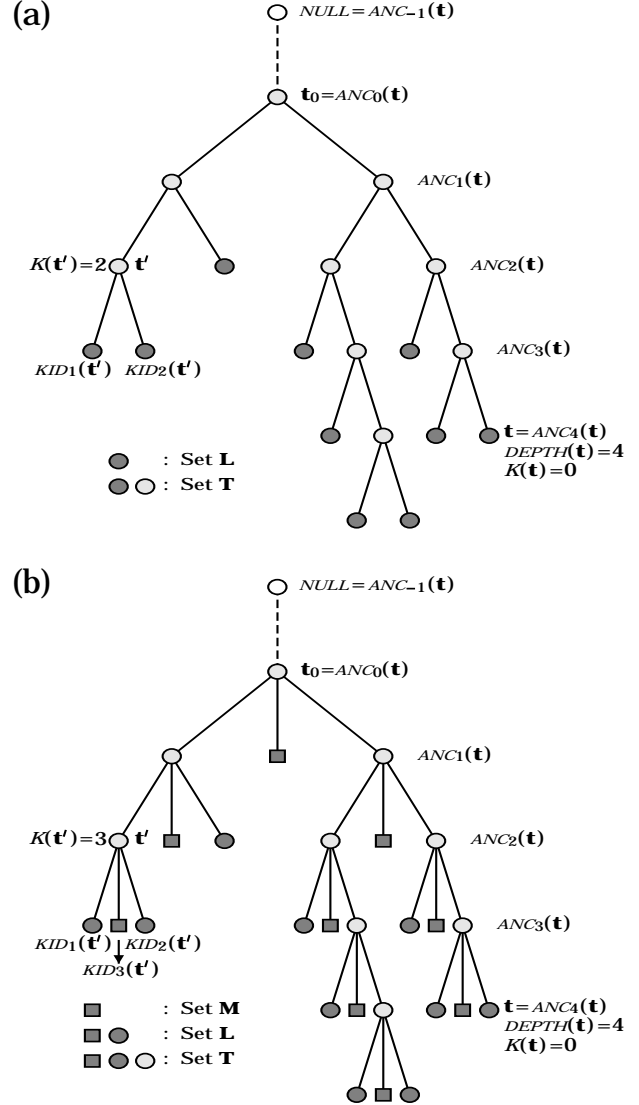


Figure 2: (a): Notation for a binary decision tree (Section 3.2). (b): Notation for a ternary decision tree (Section 3.3).

$v \in \mathcal{V}$ , immediately following any of the histories  $\mathbf{h}_{n-1}$ , clustered at node  $t$ , whereas,  $C_d(t) = \sum_{v \in \mathcal{V}} C_d(v, t) = \sum_{\mathbf{h}_{n-1} \in \Phi_t} C_d(\mathbf{h}_{n-1})$ .

Similarly to the  $n$ -gram language model, and given the equivalence classification  $\Phi$  (i.e., the tree  $\mathbf{T}$ ), probability mass functions (14), for  $t \in \mathbf{L}$ , minimize the minus log-probability on the development set,  $\mathcal{L}\mathcal{P}_d$ . The resulting  $\mathcal{L}\mathcal{P}_d$  equals the *tree entropy*, given by (Breiman et al., 1984)

$$\mathcal{H}_d(\mathbf{T}) = \sum_{t \in \mathbf{L}} f(t|t_0) \mathcal{H}_d(t), \quad (15a)$$

where

$$\mathcal{H}_d(t) = - \sum_{v \in \mathcal{V}} f(v|t) \log_2 f(v|t), \quad (15b)$$

for  $t \in \mathbf{T}$ , denotes the *node entropy* (Shannon, 1951).

In contrast to the  $n$ -gram language model, decision tree based language models offer additional freedom in designing the equivalence classification  $\Phi$ , by seeking the minimization of  $\mathcal{H}_d(\mathbf{T})$  over a set of possible trees. Decision trees can, therefore, obtain the  $n$ -gram history equivalence classes, as a special case. Clearly, minimization of (15a) over all possible trees is not feasible. Instead, a top-down and greedy algorithm is typically used for decision tree development (Breiman et al., 1984): Starting from the initial tree  $\mathbf{T}' = \{t_0\}$ , and until a stopping criterion is met, the algorithm performs the best split among the set of leaves  $\mathbf{L}'$  of the current tree  $\mathbf{T}'$ , by means of question  $\hat{\mathcal{Q}}(\hat{\mathbf{t}})$ , which is determined as

$$\hat{\mathbf{t}} = \arg \max_{\mathbf{t} \in \mathbf{L}'} \Delta \mathcal{H}_d(\hat{\mathcal{Q}}(\mathbf{t})), \quad (16a)$$

where

$$\hat{\mathcal{Q}}(\mathbf{t}) = \arg \max_{\mathcal{Q}(\mathbf{t})} \Delta \mathcal{H}_d(\mathcal{Q}(\mathbf{t})), \quad (16b)$$

and

$$\begin{aligned} \Delta \mathcal{H}_d(\mathcal{Q}(\mathbf{t})) &= f(\mathbf{t}|t_0) \mathcal{H}_d(\mathbf{t}) \\ &\quad - \sum_{i=1}^{K(\mathbf{t})} f(KID_i(\mathbf{t})|t_0) \mathcal{H}_d(KID_i(\mathbf{t})). \end{aligned} \quad (16c)$$

There are three main issues related to decision tree design (Breiman et al., 1984). (a): The choice of candidate questions  $\mathcal{Q}(\mathbf{t})$  in (16b); (b): The choice of the right size for the tree  $\mathbf{T}$ ; and (c): The estimation of  $\Pr[v|\mathbf{t}]$ , where  $v \in \mathcal{V}$  and  $\mathbf{t} \in \mathbf{L}$ .

In this work, we exclusively consider single feature questions of the type

$$\begin{aligned} \mathcal{Q}(\mathbf{t}) : \mathbf{h}_{n-1} &= (\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1}) \in \Phi_{\mathbf{t}} \\ \rightarrow \begin{cases} KID_1(\mathbf{t}), & \text{if } \omega_{-j} \in A_1 \subset \mathcal{V}, \\ KID_2(\mathbf{t}), & \text{if } \omega_{-j} \in A_2 \subseteq \mathcal{V} - A_1, \\ KID_3(\mathbf{t}), & \text{if } \omega_{-j} \in \mathcal{V} - \{A_1 \cup A_2\}, \end{cases} \end{aligned} \quad (17)$$

for all  $j=1, 2, \dots, n-1$ , where  $A_1, A_2 \neq \emptyset$  (clearly,  $A_1 \cap A_2 = \emptyset$ ). We further restrict sets  $A_1$  and  $A_2$  to satisfy, either  $A_1 \cup A_2 = \mathcal{V}$ , thus, resulting in binary splits ( $KID_3(\mathbf{t}) = \emptyset$ ), as in Section 3.2, or  $A_1 \cup A_2 = \mathcal{V}_j(\mathbf{t})$ , where

$$\begin{aligned} \mathcal{V}_j(\mathbf{t}) &= \{\omega_{-j} : \mathbf{h}_{n-1} = (\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1}) \in \Phi_{\mathbf{t}}, \\ &\quad \text{with } C_d(\mathbf{h}_{n-1}) > 0\} \subseteq \mathcal{V}, \end{aligned} \quad (18)$$

thus, resulting in ternary splits, as in Section 3.3. Notice that, in general,  $|\mathcal{V}_j(\mathbf{t})| \gg 1$ . Since there exist  $2^{|\mathcal{V}_j(\mathbf{t})|-1} - 1$  possible binary partitions of  $\mathcal{V}_j(\mathbf{t})$ , solving (16b) is, in general, computationally infeasible. Two sub-optimal solutions to this problem are discussed in Sections 3.2 and 3.3. In this paper, multiple feature questions are not considered. Notice though, that both algorithms presented in Sections 3.2 and 3.3 can be readily extended to accommodate such questions, at a significantly higher computational cost.

With regard to determining the right size of a tree, and as a result of our experiments, we have concluded that it

is beneficial to completely develop the tree using (16) until  $\Delta \mathcal{H}_d(\hat{\mathcal{Q}}(\hat{\mathbf{t}})) = 0$ . In such a case, a simple *right-to-left* tree development algorithm can be used (Potamianos and Jelinek, 1997). We then employ a suitable smoothing algorithm, in order to determine  $\Pr[v|\mathbf{t}]$ ,  $v \in \mathcal{V}$ ,  $\mathbf{t} \in \mathbf{L}$ , as we discuss in Section 3.4. Performance of the resulting language model is assessed by the minus log-probability on the test set, given by (compare to (2))

$$\mathcal{L}\mathcal{P}_t(\mathbf{T}) = - \sum_{\mathbf{t} \in \mathbf{L}} \sum_{v \in \mathcal{V}} \frac{C_t(v, \mathbf{t})}{C_t(t_0)} \log_2 \Pr[v|\mathbf{t}].$$

### 3.2. Decision tree question design by means of vocabulary binary encoding

Let us assume that we have a unique binary representation of all  $v \in \mathcal{V}$ , by means of a binary encoding of constant length  $L$ ,  $\underline{b}(v) = b_1(v) b_2(v) \dots b_L(v)$ , where  $b_i(v) \in \{0, 1\}$ , for all  $i=1, 2, \dots, L$ , and all  $v \in \mathcal{V}$ . Such an encoding is typically obtained by means of the *mutual information clustering* tree (Brown et al., 1992b; Jelinek, 1997). As an example, in Table 2, we provide the binary encoding of our Brown corpus vocabulary. Notice now that sets  $A_1 = \{v \in \mathcal{V} : b_i(v) = 0\}$  and  $A_2 = \{v \in \mathcal{V} : b_i(v) = 1\}$  satisfy  $A_1 \cup A_2 = \mathcal{V}$ , for all  $i=1, 2, \dots, L$ . Thus, the vocabulary binary encoding provides a means of partitioning set  $\mathcal{V}$  into  $L < |\mathcal{V}|$  different partitions within the totality  $2^{|\mathcal{V}|-1} - 1$  possible binary partitions of  $\mathcal{V}$ . We consider such partitions as our candidate questions of form (17) (see also (19), below).

Given the vocabulary binary encoding, any history  $\mathbf{h}_{n-1} = (\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1})$  is equivalent to  $\underline{b}(\mathbf{h}_{n-1}) = b_1(\omega_{-1}) \dots b_L(\omega_{-1}) b_1(\omega_{-2}) \dots b_L(\omega_{-2}) \dots b_1(\omega_{-n+1}) \dots b_L(\omega_{-n+1})$ , consisting of  $L \times (n-1)$  bits. Thus, we can easily build a decision tree by using as possible questions (see also (16) and (17))

$$\begin{aligned} \mathcal{Q}(\mathbf{t}, j, i) : \mathbf{h}_{n-1} &= (\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1}) \in \Phi_{\mathbf{t}} \\ \rightarrow \begin{cases} KID_1(\mathbf{t}), & \text{if } b_i(\omega_{-j}) = 0, \\ KID_2(\mathbf{t}), & \text{if } b_i(\omega_{-j}) = 1, \end{cases} \end{aligned} \quad (19)$$

for all  $j=1, 2, \dots, n-1$ , and all  $i=1, 2, \dots, L$ ; i.e.,  $\mathcal{Q}(\mathbf{t}, j, i)$  denotes the question at node  $\mathbf{t}$ , that concerns bit  $i$  of the binary encoding of  $\omega_{-j}$ . Notice that questions (19) constitute a natural choice in the case  $j=1$ , due to (16) and the fact that the maximum mutual information criterion is used in the construction of the vocabulary clustering tree (Brown et al., 1992b).

In this paper, we consider two decision tree development algorithms based on questions (19): For any tree leaf  $\mathbf{t}$ , the first algorithm searches for the best question (19), over all  $j=1, 2, \dots, n-1$ , and all  $i=1, 2, \dots, L$ , and, is, therefore, called *unrestricted*. The tree split induced by such best question is accepted, if the resulting tree entropy reduction is positive (see (16)). We refer to this algorithm as ENC.1 (Potamianos and Jelinek, 1997). The second algorithm accepts, among questions (19), the first question that results in positive tree entropy reduction, while iterating within  $j=1, 2, \dots, n-1$ , and  $i=1, 2, \dots, L$ . It is, therefore, called *restricted*. Clearly, the second algorithm forces all possible questions (19) about  $\omega_{-j}$  before any ones

$v$	$\underline{b}(v)$	$v$	$\underline{b}(v)$	$v$	$\underline{b}(v)$	$v$	$\underline{b}(v)$
!	00011110111000	7	00101111111100	O	00111111000000	h	01100000000000
#	10000000000000	8	00101111110000	P	01111010111000	i	11100000000000
\$	00011110100000	9	00101110000000	Q	01111010101000	j	01001111100000
%	00011110111110	:	00010000000000	R	01111010001110	k	01011101000000
&	00011100000000	;	00011110111110	S	01111011100000	l	01011110000000
'	00110000000000	?	00011110110000	T	01111011000000	m	01011100000000
(	00011111000000	A	00111100000000	U	00111111100000	n	01010000000000
)	00011110111100	B	01111010100000	V	01111010000110	o	11111000000000
*	00011110000000	C	01111011110000	W	01111011111000	p	01111111000000
+	00011111100000	D	01111010000000	X	01111010101100	q	01001111110000
,	00011000000000	E	00111110000000	Y	01111010011100	r	01011111000000
-	01001110000000	F	01111010111100	Z	01111010000111	s	01001100000000
.	00000000000000	G	01111010111110	'	00100000000000	t	01110000000000
0	00101000000000	H	01111010010000	a	11110000000000	u	11101000000000
1	00101100000000	I	00111000000000	b	01111110000000	v	01011101100000
2	00101111000000	J	01111010001000	c	01111100000000	w	01111000000000
3	00101111010000	K	01111010000100	d	01000000000000	x	01001111000000
4	00101111110000	L	01111010001100	e	11000000000000	y	01001000000000
5	00101111100000	M	01111010110000	f	01011000000000	z	01011101110000
6	00101111101000	N	01111010011000	g	01000100000000		

 Table 2: Binary encoding of our Brown corpus vocabulary,  $\mathcal{V}$ . Symbol “#” denotes space

about  $\omega_{-j-j'}$ , where  $j' \geq 1$ . Thus, it mimics the construction of the  $n$ -gram language model. In addition, it favors questions about the most significant bits of the vocabulary binary encoding, i.e., questions (19) with smaller values of  $i$ . We refer to this algorithm as ENC.2 (Potamianos and Jelinek, 1997).

### 3.3. Decision tree question design by means of vocabulary $K$ -means clustering

Let us again return to the problem of designing optimal questions (17). We first consider the problem in a more general setting: Given a tree node  $\mathbf{t}$  and a history position  $1 \leq j \leq n-1$ , we would like to determine a  $K$ -way partition of sub-vocabulary  $\mathcal{V}_j(\mathbf{t})$ , defined by (18), into sets  $A_1, A_2, \dots, A_K$ , where  $A_i \neq \emptyset$ ,  $\cup_{i=1}^K A_i = \mathcal{V}_j(\mathbf{t})$ , and  $A_i \cap A_{i'} = \emptyset$ , for  $i \neq i'$ , such that the  $(K+1)$ -way split of node  $\mathbf{t}$  into new nodes  $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_K$ , and a middle node  $\mathbf{t}'$  ( $\mathbf{t}' \cup_{i=1}^K \mathbf{t}_i = \mathbf{t}$ ), by means of question

$$\mathcal{Q}(\mathbf{t}; A_1, \dots, A_{K-1}) : \mathbf{h}_{n-1} = (\omega_{-1}, \dots, \omega_{-n+1}) \in \Phi_{\mathbf{t}} \rightarrow \begin{cases} \mathbf{t}_i = \text{KID}_i(\mathbf{t}), & \text{if } \omega_{-j} \in A_i, \quad i = 1, 2, \dots, K, \\ \mathbf{t}' = \text{KID}_{K+1}(\mathbf{t}), & \text{if } \omega_{-j} \in \mathcal{V} - \mathcal{V}_j(\mathbf{t}), \end{cases} \quad (20)$$

is optimal. Thus, we seek partition (see also (16))

$$\begin{aligned} & (\hat{A}_1, \hat{A}_2, \dots, \hat{A}_K) \\ & = \arg \max_{(A_1, A_2, \dots, A_K) : \cup_{i=1}^K A_i = \mathcal{V}_j(\mathbf{t}); A_i \neq \emptyset; A_i \cap A_{i'} = \emptyset} \left\{ \mathcal{H}_d(\mathbf{t}) - \sum_{i=1}^K f(\mathbf{t}_i | \mathbf{t}) \mathcal{H}_d(\mathbf{t}_i) \right\}. \end{aligned} \quad (21)$$

Let us define *atoms*, or, *atomic nodes*,  $\mathbf{v}$ , for all  $v \in \mathcal{V}_j(\mathbf{t})$  (see (18)), as

$$\mathbf{v} = \{ \mathbf{h}_{n-1} = (\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1}) \in \Phi_{\mathbf{t}} :$$

$$C_d(\mathbf{h}_{n-1}) > 0, \text{ and } \omega_{-j} = v \}.$$

Clearly, for every such atomic node, a probability mass function

$$f(\omega | \mathbf{v}) = \frac{C_d(\omega, \mathbf{v})}{C_d(\mathbf{v})}, \text{ for all } \omega \in \mathcal{V},$$

can be defined (see also (14)). Notice now that (21) is equivalent to (Chou, 1991)

$$(\hat{A}_1, \hat{A}_2, \dots, \hat{A}_K) = \quad (22)$$

$$\arg \min_{(A_1, A_2, \dots, A_K) : \cup_{i=1}^K A_i = \mathcal{V}_j(\mathbf{t}); A_i \neq \emptyset; A_i \cap A_{i'} = \emptyset} \left\{ \sum_{i=1}^K f(\mathbf{t}_i | \mathbf{t}) \mathcal{H}_d(\mathbf{t}_i) - \sum_{v \in \mathcal{V}_j(\mathbf{t})} f(\mathbf{v} | \mathbf{t}) \mathcal{H}_d(\mathbf{v}) \right\}.$$

Clearly, the argument in (22) equals (see also (15b))

$$\sum_{i=1}^K \sum_{v \in A_i} f(\mathbf{v} | \mathbf{t}) \sum_{\omega \in \mathcal{V}} f(\omega | \mathbf{v}) \log_2 \frac{f(\omega | \mathbf{v})}{f(\omega | \mathbf{t}_i)}.$$

We now re-partition node  $\mathbf{t}$  into  $\{\mathbf{t}'_1, \mathbf{t}'_2, \dots, \mathbf{t}'_K\}$ , defined by

$$\mathbf{t}'_i = \{ \mathbf{h}_{n-1} = (\omega_{-1}, \omega_{-2}, \dots, \omega_{-n+1}) \in \Phi_{\mathbf{t}} : C_d(\mathbf{h}_{n-1}) > 0, \text{ and } \omega_{-j} \in A'_i \}, \quad (23a)$$

where sets

$$A'_i = \left\{ v \in \mathcal{V}_j(\mathbf{t}) : i = \arg \min_{k \in [1, K]} \sum_{\omega \in \mathcal{V}} f(\omega | \mathbf{v}) \log_2 \frac{f(\omega | \mathbf{v})}{f(\omega | \mathbf{t}_k)} \right\}, \quad (23b)$$

for  $i = 1, 2, \dots, K$ , partition set  $\mathcal{V}_j(\mathbf{t})$ . The new partition is “better” in terms of (22), since (Chou, 1991)

$$\begin{aligned} & \sum_{i=1}^K \sum_{\mathbf{v} \in A'_i} f(\mathbf{v}|\mathbf{t}) \sum_{\omega \in \mathcal{V}} f(\omega|\mathbf{v}) \log_2 \frac{f(\omega|\mathbf{v})}{f(\omega|\mathbf{t}'_i)} \\ & \leq \sum_{i=1}^K \sum_{\mathbf{v} \in A_i} f(\mathbf{v}|\mathbf{t}) \sum_{\omega \in \mathcal{V}} f(\omega|\mathbf{v}) \log_2 \frac{f(\omega|\mathbf{v})}{f(\omega|\mathbf{t}_i)}. \end{aligned}$$

The new partition *centroids* (i.e., the probability mass functions of the resulting nodes) are given by

$$f(\omega|\mathbf{t}'_i) = \frac{\sum_{\mathbf{v} \in A'_i} f(\omega|\mathbf{v}) f(\mathbf{v}|\mathbf{t})}{\sum_{\mathbf{v} \in A'_i} f(\mathbf{v}|\mathbf{t})}, \quad (24)$$

for all  $\omega \in \mathcal{V}$ , and  $i = 1, 2, \dots, K$ . Equations (23) and (24) constitute *Chou's partitioning algorithm*, an iterative technique for improving on an initial partition  $A_1, A_2, \dots, A_K$ , of  $\mathcal{V}_j(\mathbf{t})$  on basis of (22): Until no further entropy improvement occurs by split (20), repeat the following two steps: First, compute the centroids of the current partition by means of (24) (replace  $A'_i$  by  $A_i$ ), and, second, re-partition  $\mathcal{V}_j(\mathbf{t})$  according to (23b). Notice that the algorithm results in a locally optimal partition (22) (Chou, 1991).

There exists a subtle difficulty, when repartitioning a current partition by means of (23b). It may happen that for some current node  $\mathbf{t}_k$ , atom  $\mathbf{v}$ , and  $\omega \in \mathcal{V}$ , the following holds:  $f(\omega|\mathbf{v}) > 0$  and  $f(\omega|\mathbf{t}_k) = 0$ . The *divergence* in (23b) then becomes infinite. To avoid smoothing the partition centroids (24), as is suggested by Chou (1991), we consider Chou's partitioning algorithm specialized to the *Gini index of diversity* (Breiman et al., 1984; Chou, 1991), rather than to the entropy. The Gini index of diversity is given by (compare to (15b))

$$\mathcal{G}_d(\mathbf{t}) = 1 - \sum_{v \in \mathcal{V}} f(v|\mathbf{t})^2 = \sum_{v \in \mathcal{V}} f(v|\mathbf{t}) [1 - f(v|\mathbf{t})].$$

In such a case, the *euclidean distance* between probability mass functions  $f(\omega|\mathbf{v})$  and  $f(\omega|\mathbf{t}_k)$  replaces their divergence in (23b) (Potamianos and Jelinek, 1997).

In this paper, we employ Chou's partitioning algorithm for the binary ( $K=2$ ) partitioning of set  $\mathcal{V}_j(\mathbf{t})$ . The candidate questions then become (see (17) and (20))

$$\begin{aligned} \mathcal{Q}(\mathbf{t}, j, A_1) : \mathbf{h}_{n-1} &= (\omega_{-1}, \dots, \omega_{-n+1}) \in \Phi_{\mathbf{t}} \\ &\rightarrow \begin{cases} KID_1(\mathbf{t}), & \text{if } \omega_{-j} \in A_1 \subset \mathcal{V}_j(\mathbf{t}), \\ KID_2(\mathbf{t}), & \text{if } \omega_{-j} \in A_2 = \mathcal{V}_j(\mathbf{t}) - A_1, \\ KID_3(\mathbf{t}), & \text{if } \omega_{-j} \in \mathcal{V} - \mathcal{V}_j(\mathbf{t}). \end{cases} \end{aligned} \quad (25)$$

Since Chou's partitioning algorithm is locally only optimal, we use multiple runs of the algorithm by starting from various initial partitions, constructed by partitioning  $\mathcal{V}_j(\mathbf{t})$  at random. The entropy is used as a figure of merit in comparing the resulting sub-optimal splits (Potamianos and Jelinek, 1997).

We consider three versions of Chou's partitioning algorithm for decision tree design: The first one, referred to as

*unrestricted*, or CHOU.1, at every candidate leaf, considers questions  $\mathcal{Q}(\mathbf{t}, j, \hat{A}_1)$ , where sets  $\hat{A}_1$  result from Chou's partitioning algorithm, for all  $j = 1, 2, 3, \dots, n-1$ . Thus, at every candidate leaf,  $n-1$  candidate questions (splits) are considered, and the one that maximizes the induced (by the split) tree entropy reduction is chosen (see also (16)). The second algorithm, referred to as *restricted*, or CHOU.2, at every candidate leaf, picks as the candidate question the first one that results in  $\Delta \mathcal{H}_d(\mathcal{Q}(\mathbf{t}, j, \hat{A}_1)) > 0$ , during iteration  $j = 1, 2, 3, \dots, n-1$ . It is, therefore, biased towards asking questions about  $\omega_{-j}$  before asking questions about  $\omega_{-j-j'}$ , where  $j' \geq 1$ . Finally, the third algorithm, referred to as *n-gram equivalent*, or CHOU.3, creates a tree that contains all  $n$ -gram tree nodes in it: At every candidate leaf  $\mathbf{t}$ , the algorithm picks as the candidate question the first one that satisfies either  $\Delta \mathcal{H}_d(\mathcal{Q}(\mathbf{t}, j, \hat{A}_1)) > 0$ , or  $\Delta \mathcal{H}_d(\mathcal{Q}(\mathbf{t}, j, \hat{A}_1)) = 0$  and  $|\mathcal{V}_j(\mathbf{t})| > 1$ , during iteration  $j = 1, 2, 3, \dots, n-1$ . Notice that the second condition arises, whenever  $|\mathcal{V}_j(\mathbf{t})| > 1$ , and  $f(\omega|\mathbf{v}) = f(\omega|\mathbf{v}')$ , for all  $\omega \in \mathcal{V}$ , and all  $\mathbf{v}, \mathbf{v}' \in \mathcal{V}_j(\mathbf{t})$ . In such a case, any question (25) induces zero tree entropy reduction, since  $\mathcal{H}_d(\mathbf{t}) = \mathcal{H}_d(KID_1(\mathbf{t})) = \mathcal{H}_d(KID_2(\mathbf{t})) = \mathcal{H}_d(\mathbf{v})$ , for all  $\mathbf{v} \in \mathcal{V}_j(\mathbf{t})$  (see also (15b) and (21)). We clearly have to “force” a question (25) (here, we pick a random partition of  $\mathcal{V}_j(\mathbf{t})$ ), in order to ensure that all  $n$ -gram nodes are contained in the resulting tree. Such nodes satisfy  $\mathcal{V}_j(\mathbf{t}) = 1$  before any questions about  $\omega_{-j-j'}$ , where  $j' \geq 1$ , are asked (Potamianos and Jelinek, 1997).

### 3.4. Decision tree smoothing

In Sections 3.2 and 3.3, we discussed algorithms for decision tree design, based on the development data set. Clearly, every leaf  $\mathbf{t} \in \mathbf{L} - \mathbf{M}$  can be assigned maximum likelihood estimates of probabilities (1), given by  $\text{Pr}_{\text{ML}}[v|\mathbf{t}]$ , in (14). Of course, such estimates do not generalize to unseen data well, and  $\mathcal{L}\mathcal{P}_t(\mathbf{T})$ , most certainly, becomes infinite. Therefore, smoothing techniques must be used in conjunction with (14), in order to provide the required generalization.

In Section 2, we considered a variety of smoothing techniques applied to the  $n$ -gram letter language model. We concluded that, among these techniques, the bottom-up deleted interpolation is the most robust to sparseness of training data, outperforming the rest. We, therefore, proceed to smooth our decision tree leaf probability mass functions by means of this algorithm.

Revisiting (12), we discover that decision trees exhibit certain peculiarities that make the straightforward application of the bottom-up deleted interpolation algorithm rather difficult: Decision trees are not balanced, thus, in general, tree leaves have different numbers of ancestors, some of the leaves often being of very small or very large depth. Therefore, when trying to simultaneously smooth the probability mass functions of the nodes in set  $\mathbf{L} - \mathbf{M}$  (i.e., of all tree leaves that are not middle nodes), as we do in the  $n$ -gram case, we face the obvious difficulty of how to define an equal number of distinct ancestors for each one of these leaves. In addition, the ternary decision trees, dis-

cussed in Section 3.3, pose one more challenge: At every split (25), a middle node is defined. Histories of both the held-out and test sets might end up in some of the middle nodes, which, therefore, should have a probability mass function associated with them. We choose to assign to them the probability mass function of their parent node, i.e.,

$$\Pr[v|\mathbf{t}] = \Pr[v|PAR(\mathbf{t})], \quad \text{for all } v \in \mathcal{V}, \quad \text{all } \mathbf{t} \in \mathbf{M}.$$

Practically, therefore, smoothed probability mass functions of all tree nodes are needed. Next, we propose an algorithm that addresses these issues.

Let us first define

$$DEPTH_{\min} = \min \{ DEPTH(\mathbf{t}) : \mathbf{t} \in \mathbf{L-M} \} \geq 0,$$

and

$$DEPTH_{\max} = \max \{ DEPTH(\mathbf{t}) : \mathbf{t} \in \mathbf{L-M} \} \\ \leq (n-1) (|\mathcal{V}|-1),$$

and consider some integer  $D$ :  $0 \leq D \leq DEPTH_{\max}$ . We then partition the set of tree leaves (excluding middle nodes)  $\mathbf{L-M}$  into sets

$$\mathbf{S}_1 = \{ \mathbf{t} \in \mathbf{L-M} : DEPTH(\mathbf{t}) < D \}, \\ \text{and } \mathbf{S}_2 = \{ \mathbf{t} \in \mathbf{L-M} : DEPTH(\mathbf{t}) \geq D \},$$

and the set of internal tree nodes  $\mathbf{T-L}$  into sets

$$\mathbf{S}_3 = \{ \mathbf{t} \in \mathbf{T-L} : DEPTH(\mathbf{t}) < D \}, \\ \text{and } \mathbf{S}_4 = \{ \mathbf{t} \in \mathbf{T-L} : DEPTH(\mathbf{t}) \geq D \}.$$

Clearly,  $\mathbf{S}_1 \cup \mathbf{S}_2 = \mathbf{L-M}$ ,  $\mathbf{S}_1 \cap \mathbf{S}_2 = \emptyset$ ,  $\mathbf{S}_3 \cup \mathbf{S}_4 = \mathbf{T-L}$ , and  $\mathbf{S}_3 \cap \mathbf{S}_4 = \emptyset$ . For binary trees, we need smoothed probability mass functions for all nodes in  $\mathbf{S}_1 \cup \mathbf{S}_2$ , whereas, for ternary trees, we need smoothed estimates for all tree nodes. Clearly, for every node in  $\mathbf{S}_2 \cup \mathbf{S}_4$ , we can define  $D+2$  distinct ancestor nodes of it (see (13)). We can, therefore, easily apply the bottom-up deleted interpolation algorithm on the node set  $\mathbf{S}_2 \cup \mathbf{S}_4$ . Since this constitutes the most robust smoothing algorithm, we would like set  $\mathbf{S}_2 \cup \mathbf{S}_4$  to contain as many nodes as possible. On the other hand, we would like value  $D$  to be large enough, so that enough “levels” of smoothing can be used in the deleted interpolation algorithm. In our experiments, we use value  $D = DEPTH_{\min}$ . Clearly, under this choice of  $D$ ,  $\mathbf{S}_1 = \emptyset$ .

Once the value of  $D$  is chosen, we proceed to smooth the probability mass functions of nodes  $\mathbf{t} \in \mathbf{S}_1 \cup \mathbf{S}_3$ . The depth of such nodes varies from zero to  $D-1$ , therefore, a robust deleted interpolation algorithm cannot be used. Our experimental results in Section 2 (see also Table 1) indicate that a powerful and simple smoothing algorithm, that does not require node bucketing for optimization of coefficients, is the back-off method (8) in conjunction with the discounting by half law of succession (4.2). Thus, we set  $\Pr[v|\mathbf{t}] = \Pr_{\text{BOF.2}}[v|\mathbf{t}]$ , for all  $v \in \mathcal{V}$ , and  $\mathbf{t} \in \mathbf{S}_1 \cup \mathbf{S}_3$ , where (see (4.2), (8), (13), and (14))

$$\Pr_{\text{BOF.2}}[v|ANC_{-1}(\mathbf{t})] = \frac{1}{|\mathcal{V}|},$$

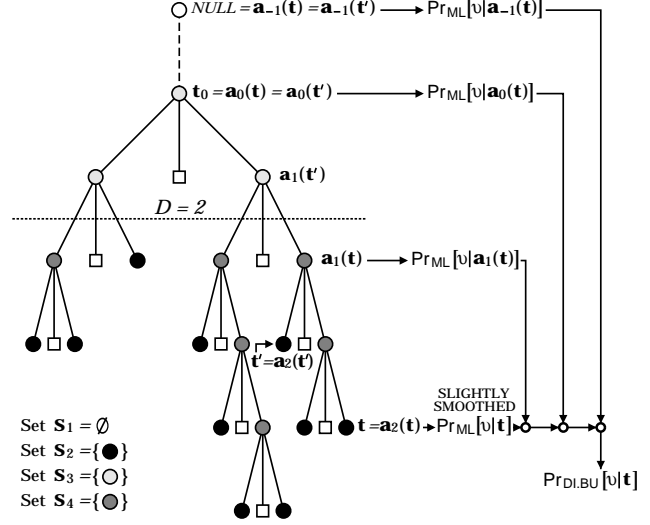


Figure 3: Tree smoothing by means of the bottom-up deleted interpolation algorithm. In this case,  $D = 2$ .

and, for  $i = 0, 1, 2, \dots, DEPTH(\mathbf{t})$ ,

$$\Pr_{\text{BOF.2}}[v|ANC_i(\mathbf{t})] = \\ \begin{cases} \frac{C_d(v, ANC_i(\mathbf{t}))}{C_d(ANC_i(\mathbf{t}))}, & \text{if } q(ANC_i(\mathbf{t})) = |\mathcal{V}|; \text{ otherwise:} \\ \frac{C_d(v, ANC_i(\mathbf{t})) - 0.5}{C_d(ANC_i(\mathbf{t}))}, & \text{if } C_d(v, ANC_i(\mathbf{t})) > 0, \\ \beta_i(\mathbf{t}) \Pr_{\text{BOF.2}}[v|ANC_{i-1}(\mathbf{t})], & \text{if } C_d(v, ANC_i(\mathbf{t})) = 0, \end{cases}$$

where, if  $q(ANC_i(\mathbf{t})) < |\mathcal{V}|$ ,

$$\beta_i(\mathbf{t}) = \frac{1 - \sum_{v: C_d(v, ANC_i(\mathbf{t})) > 0} \frac{C_d(v, ANC_i(\mathbf{t})) - 0.5}{C_d(\mathbf{t})}}{\sum_{v: C_d(v, ANC_i(\mathbf{t})) = 0} \Pr_{\text{BOF.2}}[v|ANC_{i-1}(\mathbf{t})]}.$$

For the remaining nodes  $\mathbf{t} \in \mathbf{S}_2 \cup \mathbf{S}_4$ , we use the bottom-up deleted interpolation algorithm, with  $D+1$  iterations. For this purpose, we define  $D+2$  ancestor nodes of  $\mathbf{t}$ , by (see also (12), (13b), and Fig. 3)<sup>8</sup>

$$\mathbf{a}_{-1}(\mathbf{t}) = \text{NULL}, \quad \text{and } \mathbf{a}_i(\mathbf{t}) = ANC_{\lfloor i DEPTH(\mathbf{t})/D \rfloor}(\mathbf{t}), \\ \text{for } i = 0, 1, 2, \dots, D.$$

Clearly,  $\mathbf{a}_0(\mathbf{t}) = \mathbf{t}_0$ ,  $\mathbf{a}_D(\mathbf{t}) = \mathbf{t}$ , and  $DEPTH(\mathbf{a}_{i+1}(\mathbf{t})) - DEPTH(\mathbf{a}_i(\mathbf{t})) \simeq DEPTH(\mathbf{t})/D$ , for  $i = 0, 1, \dots, D-1$ . Then, we seek to find smoothing coefficients  $\lambda^{(i)}(\mathbf{t})$ ,  $i = -1, 0, 1, \dots, D$ , that satisfy

$$0 < \lambda^{(i)}(\mathbf{t}) < 1, \quad \text{for } i = -1, 0, 1, 2, \dots, D, \\ \text{and } \sum_{i=-1}^D \lambda^{(i)}(\mathbf{t}) = 1, \quad (26a)$$

such that the probability mass function

$$\Pr_{\text{DL}}[v|\mathbf{t}] = \sum_{i=-1}^D \lambda^{(i)}(\mathbf{t}) \Pr_{\text{ML}}[v|\mathbf{a}_i(\mathbf{t})] \quad (26b)$$

<sup>8</sup>We denote  $\lfloor x \rfloor = \max \{ i : i \text{ is integer, and } i \leq x \}$ .

$n$	$n$ -gram	ENC.1	ENC.2	CHOU.1	CHOU.2	CHOU.3
1	4.380	4.380	4.380	4.380	4.380	4.380
2	3.470	3.470	3.470	3.470	3.470	3.470
3	2.850	2.850	2.850	2.850	2.850	2.850
4	2.326	2.327	2.326	2.326	2.325	2.325
5	2.007	2.012	2.010	2.008	2.006	2.004
6	1.878	1.893	1.889	1.887	1.879	1.872
7	1.831	1.861	1.855	<b>1.853</b>	1.839	1.825
8	1.811	<b>1.857</b>	1.848	1.856	1.827	1.807
9	1.801	1.858	1.844	1.866	1.820	1.800
10	<b>1.796</b>	1.862	<b>1.842</b>	1.879	<b>1.815</b>	<b>1.795</b>

Table 3: Minus log-probability on our Brown corpus test subset,  $\mathcal{L}_{\mathcal{P}_t}$ , of the  $n$ -gram and five decision tree based language models smoothed by means of bottom-up deleted interpolation, for  $n = 1, 2, \dots, 10$ . All results are measured in bits per character

generalizes to held-out data set well. Notice the similarity between (9) and (26). The bottom-up smoothing algorithm (12) can be now applied in a straightforward manner for decision tree smoothing (see also Fig. 3).

### 3.5. Experimental results

We now compare the performance of five decision tree based letter language models and the baseline  $n$ -gram letter language model on the Brown corpus. The decision trees are trained on  $n$ -gram history data,  $1 \leq n \leq 10$ , by means of the unrestricted and restricted history encoding algorithms (ENC.1 and ENC.2), discussed in Section 3.2, and the unrestricted, restricted, and  $n$ -gram equivalent Chou's partitioning algorithms (CHOU.1, CHOU.2, and CHOU.3), presented in Section 3.3. All trees are completely developed and smoothed, as discussed in Section 3.4.

In Table 3, we depict the minus log-probabilities of the six language models on the test subset of the Brown corpus. Clearly, the baseline  $n$ -gram language model defeats the first four decision tree language models. Notice that the performance of the unrestricted ENC.1 and CHOU.1 models deteriorates for  $n > 8$  and  $n > 7$ , respectively, with the CHOU.1 algorithm giving a better result. The performance of the restricted algorithms (ENC.2 and CHOU.2) further favors Chou's partitioning algorithm. Notice that only the decision trees developed by means of the CHOU.3 algorithm are equivalent to  $n$ -grams, although (having more internal nodes) they offer additional freedom in the smoothing algorithm. This translates into a slightly better performance of the CHOU.3 language model over the  $n$ -gram.

### 3.6. Remarks on decision tree based language modeling

The experimental results reported above are disappointing: The decision tree language models (with the exception of the CHOU.3 one) clearly fail to improve on the baseline

$n$ -gram language model. The cause seems to be data sparseness, coupled with the fact that the decision tree development algorithms decide on the tree splits solely on basis of seen data. Whenever, for example, a question of type (17) is asked, there might exist test set histories that correspond to set  $\mathcal{V} - \mathcal{V}_j(\mathbf{t})$ . In (19), such histories drop to the left or right children nodes depending on the highly irrelevant vocabulary encoding, whereas, in (25), they end up in the middle node, and are assigned the probability of the parent node  $\mathbf{t}$ . Both treatments are clearly inadequate, although our results indicate that the second is preferable. An  $n$ -gram language model, however, deals with unseen histories in a very natural way: It considers the largest length seen  $n'$ -gram ( $n' < n$ ), and it assigns its probability mass function to the unseen  $n$ -gram. In the case of a decision tree middle node, its probability mass function typically corresponds to a cluster of histories that apparently have less predictive capability than a specific character string ( $n'$ -gram). It is, therefore, not surprising that, when restricting the decision trees to mimic the hierarchical structure of  $n$ -grams (such as the ENC.2, CHOU.2, and CHOU.3 algorithms), the decision tree results on test data improve.

## 4. Summary and future work

In this work, we provided a comparative study of  $n$ -gram and decision tree language models, by conducting letter language modeling experiments on the Brown corpus. We first studied a variety of language model smoothing algorithms and compared their performance on the  $n$ -gram letter language model. We concluded that the bottom-up deleted interpolation algorithm constitutes the most robust smoothing technique. We then considered two algorithms for decision tree question design and concluded that the one based on  $K$ -means clustering is preferable. However, when comparing the  $n$ -gram and decision tree letter language model performance on the test set, and for large  $n$ , the former outperforms the latter. This, we believe, is due to the fact that character strings are naturally well modeled with the hierarchical structure of  $n$ -grams.

It is of great interest to consider the relative merits of the various algorithms discussed in this paper, in the case of large vocabulary word language modeling. As stated in the introduction, the data sparseness problem rapidly arises in such a case, as  $n$  increases. We believe that the employment of certain types of class language models can alleviate this problem and thus mimic the gradual appearance of data sparseness that we encountered in the letter language models. We expect our conclusions on smoothing algorithms to generalize to such cases. In addition, we believe that multiple feature questions can be beneficial for decision tree word language model design. We plan to investigate these issues in the future.

## Acknowledgements

We would like to thank Eric Ristad for his feedback on the technical report version of this paper (Potamianos and Jelinek, 1997), and Mari Ostendorf for inspiring discussions on decision tree algorithms. In addition, we would like to thank Xiaoqiang Luo for helping with Brown corpus data preparation. We also wish to thank four anonymous reviewers for their constructive comments on an earlier version of this article.

## References

- Bahl, L.R., Brown, P.F., deSouza, P.V., Mercer, R.L., 1989 . A tree-based statistical language model for natural language speech recognition. *IEEE Trans. Acoust. Speech Signal Process.* 37, 1001-1008.
- Bahl, L.R., Brown, P.F., deSouza, P.V., Mercer, R.L., Nahamoo, D., 1991. A fast algorithm for deleted interpolation. In: *Proc. Europ. Conf. Speech Comm. Tech.*, Genova, pp. 1209-1212.
- Bahl, L.R., Jelinek, F., Mercer, R.L., 1983 . A maximum likelihood approach to continuous speech recognition. *IEEE Trans. Pattern Anal. Machine Intell.* 5, 179-190.
- Bell, T.C., Cleary, J.G., Witten, I.H., 1990 . *Text Compression*. Prentice Hall, Englewood Cliffs.
- Betz, M., Hild, H., 1995 . Language models for a spelled letter recognizer. In: *Proc. Internat. Conf. Acoust. Speech Signal Process.*, Detroit, pp. 856-859.
- Breiman, L., Friedman, J.H., Olshen, R.A., Stone, C.J., 1984. *Classification and Regression Trees*. Chapman and Hall, New York.
- Brown, P.F., Della Pietra, S.A., Della Pietra, V.J., Lai, J.C., Mercer, R.L., 1992a. An estimate of an upper bound for the entropy of English. *Comput. Linguistics* 18, 31-40.
- Brown, P.F., Della Pietra, V.J., deSouza, P.V., Lai, J.C., Mercer, R.L., 1992b. Class-based  $n$ -gram models of natural language. *Comput. Linguistics* 18, 467-479.
- Chen, S.F., Goodman, J., 1996. An empirical study of smoothing techniques for language modeling. In: *Proc. 34th Ann. Meeting Assoc. Comput. Linguistics*, Santa Cruz, pp. 310-318.
- Chou, P.A., 1991 . Optimal partitioning for classification and regression trees. *IEEE Trans. Pattern Anal. Machine Intell.* 13, 340-354.
- Efron, B., 1982 . *The Jackknife, the Bootstrap and Other Resampling Plans*. CBMS-NSF Reg. Conf. Ser. Appl. Math., Society for Industrial and Applied Mathematics, Philadelphia.
- Good, I.J., 1953 . The population frequencies of species and the estimation of population parameters. *Biometrika* 40, 237-264.
- Jelinek, F., 1997 . *Statistical Methods for Speech Recognition*. MIT Press, Cambridge.
- Jelinek, F., Mercer, R.L., 1980 . Interpolated estimation of Markov source parameters from sparse data. In: *Gel-sema, E.S., Kanal, L.N. (Eds.), Pattern Recognition in Practice*, North Holland, Amsterdam, pp. 381-397.
- Katz, S.M., 1987 . Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE Trans. Acoust. Speech Signal Process.* 35, 400-401.
- Kucera, H., Francis, W., 1967 . *Computational Analysis of Present-Day American English*. Brown University Press, Providence.
- Ney, H., Essen, U., Kneser, R., 1995 . On the estimation of 'small' probabilities by leaving-one-out. *IEEE Trans. Pattern Anal. Machine Intell.* 17, 1202-1212.
- Potamianos, G., Jelinek, F., 1997 . A study of  $n$ -gram and decision tree letter language modeling methods. *CLSP Research Notes 13*, Center for Language and Speech Processing, The Johns Hopkins University, Baltimore.
- Press, W.H., Flannery, B.P., Teukolsky, S.A., Vetterling, W. T., 1988. *Numerical Recipes in C. The Art of Scientific Computing*. Cambridge University Press, Cambridge.
- Ristad, E.S., 1995 . A natural law of succession. *Research Report CS-TR-495-95*, Department of Computer Science, Princeton University, Princeton.
- Ristad, E.S., Thomas, R.G., 1995 . New techniques for context modeling. In: *Proc. 33rd Ann. Meeting Assoc. Comput. Linguistics*, Cambridge, pp. 220-227.
- Ristad, E.S., Thomas, R.G., 1997a . Hierarchical non-emitting Markov models. In: *Proc. 35th Ann. Meeting Assoc. Comput. Linguistics*, Madrid, pp. 381-385.
- Ristad, E.S., Thomas, R.G., 1997b . Hierarchical non-emitting Markov models. *Research Report CS-TR-544-97*, Department of Computer Science, Princeton University, Princeton.
- Shannon, C.E., 1951 . Prediction and entropy of printed English. *Bell Syst. Tech. J.* 30, 50-64.