

論文 / 著書情報
Article / Book Information

論題(和文)	音響モデル学習のための相対エントロピーを用いた学習文選択手法
Title(English)	
著者(和文)	村上博子, 篠田浩一, 古井貞熙
Authors(English)	Hiroko Murakami, Koichi Shinoda, SADAOKI FURUI
出典(和文)	日本音響学会2011年春季講演論文集, Vol. , No. 1-5-7, pp. 17-20
Citation(English)	, Vol. , No. 1-5-7, pp. 17-20
発行日 / Pub. date	2011, 3

音響モデル学習のための相対エントロピーを用いた 学習文選択手法*

村上博子, 篠田浩一, 古井貞熙 (東工大)

1 はじめに

音声認識システムにおける大規模音声データベースの構築には, 大きく分けて次の 2 つの方法が考えられる. 録音した音声データを人手で書き起こす方法と, 予め用意されたテキストを多数の人が読み上げる方法である. どちらの方法でも, その構築にかかるコストは大きい.

本稿では, まず最初に後者について検討する. 音響モデルの認識性能の向上に対し, より効果の高い文をデザインすることで, 発声に用いる文数を削減できると期待できる. 従来手法としては, 予め多様な認識単位が含まれる音素バランス文セットを作成する手法が主流である [1][2][3]. しかし, 音響モデルによる認識率が低いことが予想される認識単位を, 認識率が充分高いことが予想される認識単位よりも多く出現させた方が, より認識性能の高い音響モデルを構築できる可能性がある.

そこで, 認識率が低い認識単位が多く含まれる文を学習文候補からより多く選択することで, 従来より少ない学習文数で同等の認識性能を持つ音響モデルを学習する手法を提案する. 提案手法では, 相対エントロピーを用いて, 認識単位の誤認識個数の分布と, 選択した文集合の認識単位の出現頻度の分布の類似度が最も高くなる文集合を選択する. 相対エントロピーの計算において, 近似を用いることで, 計算時間を削減した.

また, 本稿では自然発話の音声データベースの構築にも対応するため, 半教師付き学習による学習文選択の検討も行う.

2 学習アルゴリズム

提案手法では, データの増加に伴い, 認識単位を monophone, diphone, triphone と切り替えて学習文選択を行う. これは, 認識単位の数が多いと, 少ない学習データで構築された音響モデルの認識率を正確に予測することが難しいためである.

提案する学習アルゴリズムの概略を Fig. 1 に

示す. まず, 準備段階として, 選択に用いる学習候補文各々において, monophone, diphone, triphone の 3 つの認識単位について各認識単位の出現数の分布である頻度分布を求めておく. また, 初期発声データを用意し, その半分を選択済み学習文セット, 残りを認識文セットとして用いる. 1 つの認識単位に対する文選択は以下の 6 ステップで実行される.

1. 選択済み学習文セットを用いて monophone の音響モデルを構築し, 認識文セットを音素認識する.
 2. 1. の認識結果から, 認識単位の誤認識個数の分布である誤り分布を求める.
 3. 学習文候補それぞれについて, その文を加えた選択済みの文集合の頻度分布と 2. で求めた誤り分布との間の分布間距離を計算する. 分布間距離として, カルバック・ライブラー距離 (Kullback-Leibler divergence, KLD) [4] を用いる. 候補の中で最も距離が小さい文を選択し, その値を D_{select} とする.
 4. D_{select} と, 選択直前の KLD 値 D を比較する. D_{select} の方が小さければ, $D = D_{\text{select}}$ とし, 3. に戻る. そうでなければ選択を打ち切る.
 5. 話者に選択した文の発声を促す.
 6. 選択した文を学習文候補から取り除き, 選択済み学習文セットに追加し, 1. に戻る
3. で用いる分布間距離の詳細については 4 章で述べる.

提案手法は, 上記の 1. から 6. までのステップを, 最初は monophone を用いて行い, 残りの候補文に対し, 2 回目は diphone, 3 回目は triphone を用いて行う. 最後に, 選択された全てのデータを用いて triphone で音響モデルを作成し, 単語単位の認識を行う.

* A relative-entropy-based sentence selection approach for acoustic model training . by Hiroko Murakami, Koichi Shinoda, and Sadaoki Furui (Tokyo Institute of Technology)

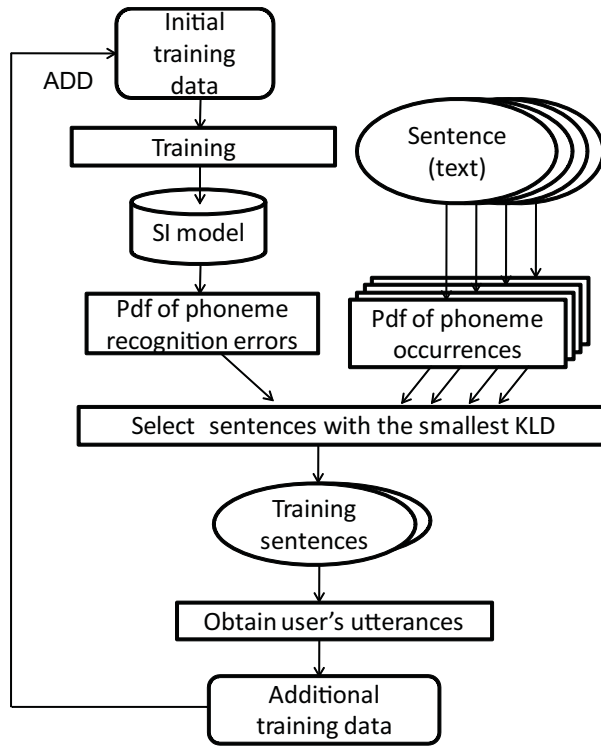


Fig. 1 Flow of the proposed method .

3 誤り分布と頻度分布

本章では，誤り分布と頻度分布の導出について述べる．まず，誤り分布を定義する．2章のステップ1で得られる音素認識結果から，各認識単位の誤認識個数を得る．別の単位に誤って認識した場合だけでなく，別の単位をその単位と誤って認識した場合も誤認識個数として数える．分布作成に用いる認識単位は，全ての認識単位から，出現数が多い認識単位を選ぶ．分布を求める際に利用する認識単位の集合を U とし，認識単位 u の誤認識個数を $r(u)$ とすると，全認識単位にわたる誤り分布 $p(u)$ は以下ようになる．

$$p(u) = \frac{r(u)}{\sum_{u \in U} r(u)} \quad (1)$$

次に，頻度分布を定義する． $s(u)$ をある候補文に含まれる認識単位 u の出現回数とすると，その文の頻度分布 $q(u)$ は以下ようになる．

$$q(u) = \frac{s(u)}{\sum_{u \in U} s(u)} \quad (2)$$

4 分布間距離

提案手法では，初期発声データを含む選択済みの文の集合に追加すると誤り分布と頻度分布の間の KLD が減少する文を，候補から選択する．

選択済みの文までの頻度分布を $q'(u)$ ，選択済みの文の総数を N とし，誤り分布 $p(u)$ と文集合の頻度分布の間の KLD 値 D を以下のように定義する．

$$D = \sum_{u \in U} p(u) \log \frac{p(u)}{q'(u)} \quad (3)$$

そこに，頻度分布 $q(u)$ をもつ文が文集合に加えられたときの KLD 値 D^+ は以下ようになる．

$$D^+ = \sum_{u \in U} p(u) \log \frac{p(u)}{(Nq'(u) + q(u))/(N+1)} \quad (4)$$

diphone, triphone では分布の単位数が多いため，新たな文が選択される度に KLD 値を直接計算すると，計算量が多くなる．ここで，直前の KLD 値との差分のみを計算することで，計算量を削減できる．差分は以下の式で表すことができる．

$$\Delta = D^+ - D = \log \left(1 + \frac{1}{N} \right) - \sum_{u \in U} p(u) \log \left(1 + \frac{q(u)}{Nq'(u)} \right) \quad (5)$$

$\Delta \geq 0$ となるとき，選択を終了する．

各文の頻度分布 $q(u)$ の各認識単位の値のほとんどが 0 となることに着目する．差分のみを計算することで， $q(u)$ の各認識単位の値が 0 となる時， \log の計算を省略することができる．さらに，式 (5) において，テーラー展開式を用いて近似を行うことで，さらに計算量を削減することができる．テーラー展開の第 1 項，及び第 2 項まで用いて近似した差分の式 Δ^1 , Δ^2 は以下ようになる．

$$\Delta^1 \approx \frac{1}{N} \sum_{u \in U} p(u) \left(1 - \frac{q(u)}{q'(u)} \right) \quad (6)$$

$$\Delta^2 \approx \Delta^1 - \frac{1}{2N^2} \sum_{u \in U} p(u) \left(1 - \left(\frac{q(u)}{q'(u)} \right)^2 \right) \quad (7)$$

5 半教師付き学習

本章では，半教師付き学習による文選択手法のアルゴリズムを述べる．基本的には，2.2 節で述べた教師付きの手法と同じアルゴリズムで選択を行う．異なる点は，学習文候補として書き起こしなしの音声データを用いるため，候補文を認識して，仮の書き起こしを得る必要があるという点である．仮の書き起こしの精度はできる

だけ高いことが望ましい．そのため，認識精度の低くなる音素認識結果ではなく，triphone の音響モデルを用いた連続音声認識（単語単位）の結果の音素列を仮の書き起こしとして用いる．

まず，書き起こしありの音声データ全てを用いて，triphone で音響モデルを学習する．そのモデルを用いて，書き起こしなしの音声データを連続音声認識（単語単位）し，その認識結果を仮の書き起こしとして用いる．後は，教師付きの文選択手法と同じアルゴリズムで選択を行い，選択した音声データを実際に書き起こし，選択済み学習文セットとして用いる．

6 実験

6.1 実験条件

データベースとして，日本語話し言葉コーパス (CSJ) [5] における男性話者による学会講演音声を用いた．全データのうち，198,807 発話 (666 話者，152 時間) を学習データとし，2,328 発話 (10 話者，1.95 時間) をテストセットとした．

音声認識に使う特徴量は MFCC12 次元とパワー，及び，その一次微分と二次微分の計 39 次元を用いた．音響モデルは 16 混合 3,000 状態 triphone HMM を用いた．実験には HTK [6] を用いた．

全学習データからランダムに 13,028 発話 (10 時間) を選択し，半分を選択済み学習文セット，残りを認識文セットとして用いる．残りの学習データ 185,779 発話 (142 時間) は学習文選択の候補文として用いる．比較実験として，候補文から学習文をランダムに選択するランダム選択と比較した．言語モデルは全学習データを用いて学習したものをを用いた．diphone は，前の音素からの遷移を考慮した左 diphone を用いた．誤り分布，頻度分布の作成に用いる認識単位は，全ての書き起こしテキストの中に 10,000 個以上出現するものをを用いた．

4 章で述べた近似は，認識単位が diphone，及び triphone のときに行った．近似を行わず log の計算をした場合 (Log) と，テーラー展開の 1 項目 (近似 1)，及び 2 項目 (近似 2) まで近似した場合を比較した．

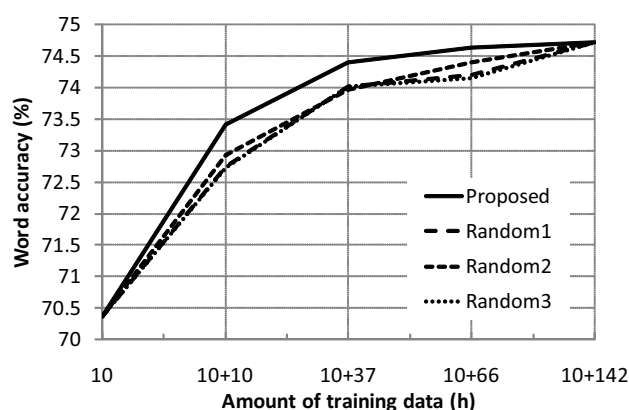


Fig. 2 Comparison of the proposed method with random selection .

6.2 実験結果

6.2.1 ランダム選択との結果比較

Fig. 2 に提案手法とランダム選択の認識結果を示す．ランダム選択では，提案手法の各認識単位において選択された文と同じ時間分の学習文を候補からランダムに選択する．ランダム選択は 3 回行った．提案手法はランダム選択と比べて良い結果となった．候補文全ての 152 時間を学習に用いて到達できる単語正解精度 74.7% を，提案手法は 50% の 76 時間の学習で達成できた．また，monophone による選択終了時 (10h)，diphone による選択終了時 (37h) においてもランダム選択より良い結果を得ることができた．

6.2.2 KLD 値の変化

Fig. 3 に，選択文数の増加に伴う，誤り分布と文集合の頻度分布の間の KLD 値を示す．認識単位をより詳細な (種類数の多い) 認識単位に変更すると，KLD 値の減少の割合は小さくなり，選択される文数が多くなる．これは，選択回数を重ねると，選択済みの文数が多くなり，式 (4) における $Nq'(u)$ が大きくなるためである．KLD の最小値が認識単位を変更するごとに大きくなっていくのは，認識単位数が多くなると頻度分布を誤り分布に正確に近づけるのが難しくなるためである．

6.2.3 近似による結果比較

Table. 1 に近似を行ったときの提案手法の実験結果を示す．認識精度にばらつきはあるが，いずれもランダム選択よりも高い値となっている．計算時間は，近似 1 は diphone で 56.0%，triphone で 39.3%，近似 2 は diphone で 48.9%，triphone で 27.9% の削減になっている．近似では，選択さ

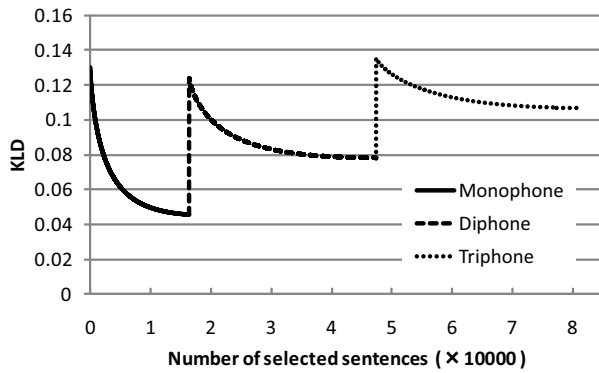


Fig. 3 Change of KLD values according to the number of selected sentences .

Table 1 Comparison of the proposed methods and random selection . Log indicates results without approximation . Ap1 and Ap2 indicate results using approximation by Taylor expansion . The table shows recognition accuracy , and time needed for sentence selection .

	Diphone			Triphone		
	Log	Ap1	Ap2	Log	Ap1	Ap2
Acc(%)	74.3	74.2	74.4	74.6	74.7	74.5
Time(h)	4.0	1.8	2.0	5.3	3.2	3.8

れる順番に違いはあるが、近似を行わない場合とほとんど同じ文が選択されている。各認識単位における選択終了後の、誤り分布と文集合の頻度分布の間の KLD 値は、近似の有無、近似方法の違いによる大きな違いはなかった。

6.2.4 半教師付き選択の実験結果

Fig. 4 に半教師付きの文選択手法の認識結果を示す。提案手法は、triphone における選択終了時に、ランダム選択の平均と比べ 0.1 ポイント良い結果となった。残念ながら、教師付きの文選択手法と比べると、ランダム選択との違いがほとんどなくなった。これは、誤りが含まれる認識結果文を文選択の際のラベルとして用いたため、誤り分布と最も近い文を選択できなかったことが原因と考えられる。

7 まとめ

音響モデル学習のための、相対エントロピーを用いた学習文選択手法を提案した。提案手法を日本語話し言葉コーパスの音声データを使い評価し、ランダムな学習文選択より良い結果を得た。

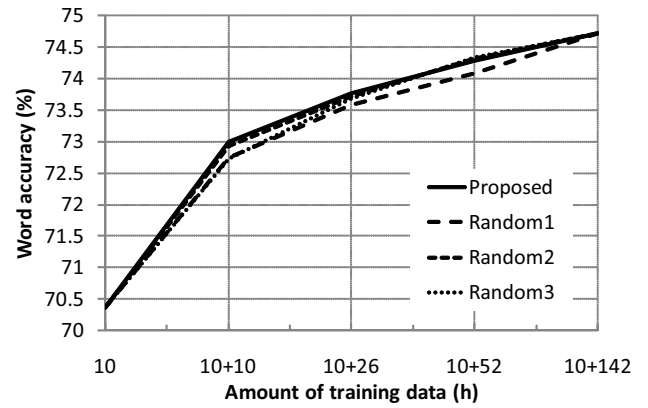


Fig. 4 Comparison of the semi-supervised selection method with random selection.

候補文全てを学習に用いて達成される 74.7% の単語正解精度に、提案手法は半分の時間の学習で到達することができた。また、相対エントロピーの計算において、近似を用いることにより、計算時間を diphone で 56.0%、triphone で 39.3% 削減した。半教師付きの文選択の実験も行ったが、教師付きの文選択ほど良い結果は得られなかった。

今後の課題として、半教師付きの文選択の改良や、学習データの選択単位を文単位から変更することによる、さらなる計算時間の削減が挙げられる。

参考文献

- [1] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, "ATR Japanese speech database as a tool of speech recognition and synthesis", Speech Communication, vol. 9, pp. 357-363, 1990 .
- [2] K. Maekawa, and K. Hanae, "Design of a spontaneous speech corpus for Japanese," Proc. the International Symposium: Toward the Realization of Spontaneous Speech Engineering, pp. 70-77, 2000 .
- [3] 磯健一, 渡辺隆夫, 桑原尚夫, "音声データベース用文セットの設計," 日本音響学会講演論文集, 2-2-19, pp.89-90, 1988 .
- [4] S. Kullback, and R. A. Leibler, J. B. MacQueen, "On information and sufficiency," Annals of Mathematical Statistics, vol. 22, no. 1, pp. 79-86, 1951.
- [5] K. Maekawa, H. Koiso, S. Furui, and H. Isahara, "Spontaneous speech corpus of Japanese," Proc. LREC, vol. 2, pp.947-952, 2000 .
- [6] Hidden Markov Model Toolkit (HTK), <http://htk.eng.cam.ac.uk/>