

ニュース音声の話題決定のためのワードスポッティングの評価

加藤 大知 山下 洋一

立命館大学理工学部情報学科

〒525-8577 滋賀県草津市野路東1-1-1

E-mail : taichi@slp.cs.ritsumeai.ac.jp, yama@cs.ritsumeai.ac.jp

概要：ニュース音声などの連続音声の話題を自動決定するために、ワードスポッティングを用いてキーワードを抽出し、抽出されたキーワード列に基づいて話題を決定する手法を検討する。本報告では、話題決定に用いるキーワードセットの決定法を、ワードスポッティングにおける検出率/湧き出し誤り率および話題決定率の観点から評価する。キーワードの選出尺度として、話題との関連性を示す尺度として相互情報量、負の値を持つ X^2 値、キーワードの検出のしやすさの尺度として湧き出し誤りの起こり難さ (DF 値) を用い、これらを組み合わせることによりキーワードセットを選出した。

キーワード：ワードスポッティング、湧き出し誤り、相互情報量、負の値を持つ X^2 値

Evaluation of Word Spotting for Topic Identification of News Speech.

Taichi Kato, Yoichi Yamashita.

Department of Computer Science, Ritsumeikan University

1-1-1 Noji-Higashi, Kusatsu-shi, Shiga, 525-8557 Japan

Abstract : This paper describes a method for automatically identifying topic for speech materials such as news speech. Topic identification is based on a sequence of keywords extracted by a word spotting technique. We used three types of measures of keyword selection, mutual information, modified X^2 score, and detection facility. The first two measures indicate degree of relevance to topic, and the last one estimates the number of false alarms of a word. Combination of these measures selected several keyword sets. The keyword sets were evaluated from viewpoints of performance of the keyword detection and the topic identification.

Keyword : word spotting, false alarm, mutual information, X^2

1 はじめに

現在、社会の情報化が進み、様々な情報が大量に蓄積されるようになってきた。蓄積される情報が増えれば増えるほど、必要な情報だけを効率よく取り出すことが難しくなり、効率の良い検索手法や、蓄積されたデータのインデキシングが必要となってくる。例えば、テレビなどのニュース音声のデータベースから必要な記事を検索するときには、記事の話題分類が有効なインデックスとなる。音声メディアに対する検索を容易にするために、ニュースなどの音声データの話題（ジャンル）を自動的に決定する方法が、近年広く研究されている [1] [2] [3]。

人が話の内容（話題）を決定するとき話に出てくる単語が重要な情報となっている。連続音声の話題を決定する時にも、音声中より話題との関連性が高い語（以下、キーワード）を抽出し、その情報を用い話題決定を行うことは、有効な手段であると考えられる。

キーワードの抽出方法は、音声中からワードスポッティングによって直接抽出する方法と連続音声認識を行い、認識結果から必要なキーワードを抽出する方法などがこれまでに報告されている。ワードスポッティングによるキーワード抽出は、連続音声認識に比べ、言語モデルと対象データのずれ、話題決定に有効な新語の扱い（未知語）、いいよどみへの対応の問題がないことから、本研究ではワードスポッティングによるキーワード抽出を行っている。

ワードスポッティングに基づいて話題決定を行う場合には、話題への関連性だけでなく検出のしやすさも考慮して検出すべきキーワードセットを決定しなければならない。著者らは、これまでに音声認識シミュレーションによってワードスポッティングによる単語の検出容易さを見積もる手法を検討してきた [4] [5]。本研究では、検出のしやすさも考慮したいくつかのキーワード選択手法を試み、ワードスポッティング性能話題決定率の観点から評価する。

2 ワードスポッティングに基づいたニュース音声の話題決定

図1に示すように、1) 予めキーワードセットを選択しておき、2) ニュース音声からキー

ワードの検出を行い、3) 検出されたキーワード列から各話題の話題らしさを算出して、最も値の高い話題を取り出し、話題決定を行う。検出されたキーワード列 W_i に対して式(1)を用いて、話題 T_j の話題らしさ $S(T_j)$ を算出する。

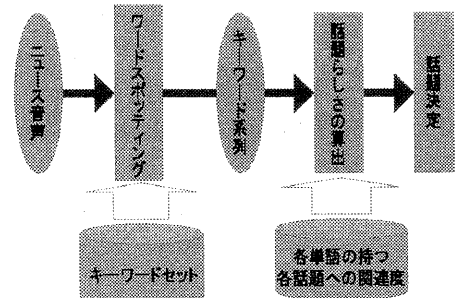


図1 話題決定の処理手順

$$S(T_j) = \sum_i (\log P(w_i | T_j) + c * S(w_i)) \quad (1)$$

c : $S(w_i)$ の重み

$P(w_i | T_j)$: 話題 T_j における単語 w_i の出現確率

$S(w_i)$: 音声中で単語 w_i が検出されたときの尤度

3 キーワード選出の尺度

キーワードの良さとして

- 1) 相互情報量
- 2) 負の値を持つ x^2 値
- 3) 湧き出し誤りの起こり難さ (DF 値)

の3つを考え、これらを組み合わせて利用することによりキーワードを選択する。キーワードの良さを示す尺度として、下記の1~6の手法について評価を行う。

手法 1 - 1	相互情報量
手法 1 - 1	相互情報量 × DF 値
手法 2 - 1	x^2 値
手法 2 - 2	x^2 値 × DF 値
手法 2 - 3	x^2 値 + ω × DF 値 ($\omega=10$)
手法 2 - 4	x^2 値 + ω × DF 値 ($\omega=50$)

まず、毎日新聞 CD-ROM 45ヶ月分 (1991年~1994年9月) に出現する名詞で、出現数上位 20, 000 単語を選出する。選出された 20, 000 単語に対し上記の手法 1 から 6 の方法でスコアを算出し、スコアの

高い、4音素以上の単語を各3,000単語ずつ選択し使用する。

3.1 相互情報量

単語 w と話題 T の相互情報量を考える。 T を記事の話題を表す確率変数、 W を記事中のある単語を表す確率変数とすると、相互情報量 $I(T;W)$ は、

$$\begin{aligned} I(T;W) &= I(W;T) \\ &= H(W) - H(W|T) \\ &= \sum_w [-p(w) \log_2 P(w) \\ &\quad + \sum_T p(T) p(w|T) \log_2 p(w|T)] \quad (2) \end{aligned}$$

となり、様々な単語から寄与を合計したものとなる。従って、式(2)の [] 内の値

$$\begin{aligned} G(w) &= -p(w) \log_2 p(w) \\ &\quad + \sum_T p(T) p(w|T) \log_2 p(w|T) \end{aligned}$$

が大きい単語 w ほど特定の話題に多くの情報を持っていると考える [1]。

3.2 負の値を持つ X^2 値

式(3)に示す負の値を持つ X^2 値は、分野における単語の偏りを示す指標として用いることができる。単語 w の出現確率は全分野に対して等しいという仮説を設定し、この仮定に基づいて単語 w の分野における予測頻度 m を計算する。また各単語 w について分野 t における頻度 x を求める。もし負の値を持つ X^2 値が十分に大きな値になれば特定の分野に偏って表れる単語ということになり、分野識別に有効な単語とみなせる [2] [3]。

$$m_{ij} = \frac{\sum_{j=1}^n x_{ij}}{\sum_{i=1}^m \sum_{j=1}^n x_{ij}} \quad (3)$$

$$X_{ij} = \frac{(x_{ij} - m_{ij}) | x_{ij} - m_{ij} |}{m_{ij}} \quad (4)$$

m : 異なり単語数

n : 分野数

x_{ij} : 単語 w_i の分野 t_j における頻度

m_{ij} : 単語 w_i の分野 t_j における予測頻度

3.3 湧き出し誤りの起こり難さ (DF値)

ワードスポッティングでは検出すべき単語が長いほど検出が容易で誤検出も少なくなる。しかし、検出の容易さは単語の音素数(モーラ数)だけでなく単語を構成する音の並びにも依存すると考えられる。そこで、単語認識シミュレーションを行い、湧き出し易さを見積もることを考える[4][5]。

ワードスポッティングで検出すべき単語 w の音素長を L 、 w を構成する音素の並びを pw_1, pw_2, \dots, pw_L とする。文における音素列の生成を n 重マルコフ過程によってモデル化し、シミュレーションによって長さ M の K 個の音素列 s_k ($k=1, \dots, K$) を生成する。 s_k を構成する音素の並びを $ps_{k1}, ps_{k2}, \dots, ps_{kM}$ とする。音素 p_i が音素 p_j に認識される確率、音素 p_i が脱落する確率、音素の挿入があるときに挿入される音素が p_j である確率をそれぞれ、 $S(p_j | p_i)$ 、 $D(p_i)$ 、 $I(p_i)$ で与えられたとする。ただし N を音素の種類としたとき

$$\sum_{j=1}^N S(p_j | p_i) + D(p_i) = 1, \sum_{j=1}^N I(p_j) = 1 \quad (5)$$

である。ここで s_k の部分音素列と単語 w との様々な対応関係を考える。その中で、 s_k の部分音素列が単語 w と認識される最大の確率 MP_k を終点フリーの動的計画法を用いて求める。 s_k の部分音素列としては、長さが $L-D$ か $L+D$ の範囲を考える。したがって、 $M, L+D$ でなければならない。 MP_k は、

$$1) \quad g(1,1) = S(pw_1 | ps_{k1}),$$

$$g(1, j) = 0 \quad (j=2, \dots, L)$$

$$2) \quad i=2, \dots, L+D \text{ に対して}$$

3),4),5)を繰り返す。

$$3) \quad j1 = \max(1, i-D), j2 = \min(L, i+D)$$

$$4) \quad j = j1, \dots, j2 \text{ に対して } 5) \text{ を繰り返す。}$$

$$5) \quad g(i, j) = \max \begin{cases} g(i-1, j-1) \times S(pw_j | ps_i) \\ g(i-1, j) \times D(ps_i) \\ g(i, j-1) \times I(pw_j) \end{cases} \quad (6)$$

$$6) \quad MP_k = \max(g(L-D, L), \dots, g(L+D, L))$$

で与えられる。ここで D は非線形伸縮の程度を制御するパラメータである。さらに湧き出し誤りの起こりやすさ FP は、K 個の音素列について MP_k の和をとったものとして、

$$FP = \sum_{k=1}^K MP_k \quad (7)$$

となる。さらに、湧き出しの起こり難さを、

$$DF = -\log FP \quad (8)$$

とする。

4 評価実験

ワースポットリングによって検出されたキーワードの持つ各話題への寄与度を用いて各話題の話題らしさを式(1)を用いて算出する。話題決定に用いる話題は、“国際”、“経済”、“家庭”、“文化”、“科学”、“芸能”、“スポーツ”、“社会”の8話題である。

4.1 学習、評価データ

各キーワードの各話題に対する寄与度の学習データとして、「毎日新聞 CD-ROM」45ヶ月分('91年1月~'94年9月)を用い、評価データには NHK 1時のニュース22日分78記事('98年8月収録、総計1.83時間、1記事平均約100秒、名詞単語の種類は2486単語)を使用した。78記事の内訳は、社会53記事、国際16記事、経済7記事、スポーツ2記事であった。

4.2 ワードスポッティング

まず、キーワードセットを比較する。表1に各キーワードセットの平均音素長、出現キーワードの平均音素長、出現のべ回数平均音素長、出現キーワードの種類、1話題に出現

する平均キーワード数、評価データ中の単語におけるキーワードの割合(カバー率)を示す。

次に文献[4]に述べられている手法でワードスポッティングを行う。図2にしきい値を変えた時の各手法でのキーワード検出率を示す。ワードスポッティングのしきい値は、-1~-20で実験を行った。図2より DF 値を考慮した手法2-2のキーワードセットが湧き出し誤りの減少、より高い検出率を得ていることがわかる。手法1の中では、手法1-1よりも手法1-2によるキーワードセットの方が検出率が低かった。手法1-2は DF 値を考慮したキーワード選択法であり、比較的長い音素長のキーワードが多く選択されていることが期待される。しかし、3000キーワードの平均音素長では手法1-1より手法1-2の方が長いものの、実際に音声データに出現したキーワードの平均音素長では、表1に示すように手法1-2の方が短くなっている。このため手法1-2によるキーワードセットのスポッティング性能が低くなったと思われる。

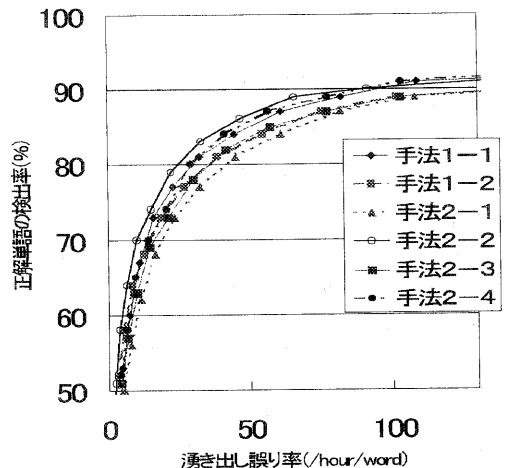


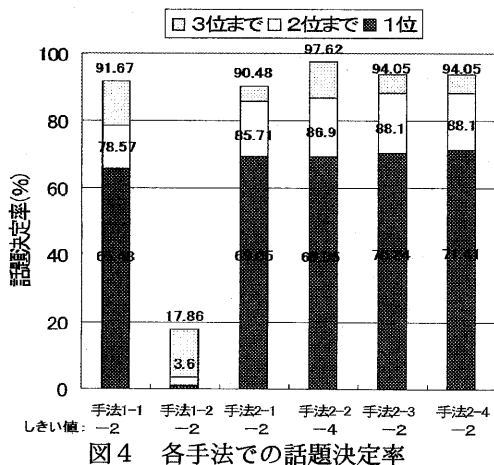
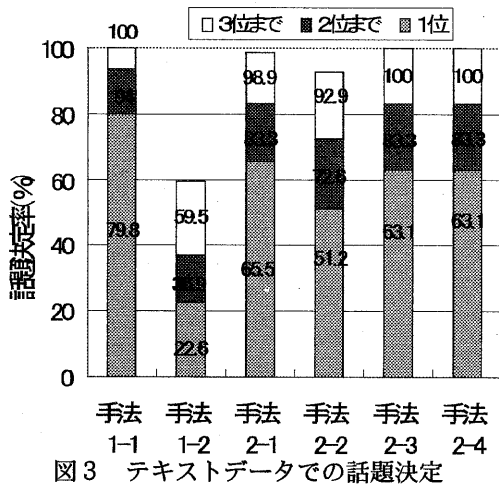
図2 正解単語の検出率

表1 評価テキストデータにおけるキーワードセットの比較

	手法1-1	手法1-2	手法2-1	手法2-2	手法2-3	手法2-4
キーワードセットの平均音素長	6.48	6.84	6.33	7.57	6.61	7.08
出現キーワードの平均音素長	6.04	6.02	5.79	6.36	6.00	6.08
キーワードの出現のべ回数	3326回	989回	3645回	1811回	3503回	3435回
出現キーワードの種類	710単語	299単語	756単語	404単語	711単語	687単語
1話題に出現する平均キーワード数	8.5単語	3.6単語	9単語	4.8単語	8.5単語	8.2単語
カバー率	29%	12%	30%	16%	29%	28%

4.3 話題決定

まず、テキストデータを対象としてキーワード列に基づく話題決定を行った。これは、ワードスポッティングを行ったときに検出率が100%で湧き出し誤りが1つも出なかった状態と同じである。その結果を図3に示す。DF値を考慮した手法1-2、2-2、3、4がDF値を考慮しない手法1-1、2-1より話題決定率がやや低い結果となった。この理由として、DF値を考慮することによりキーワードと話題との関連性が下がってしまったことが考えられる。表1でのカバー率からわかるように、手法1-2が極端に話題決定率が低いのは、評価データに出現するキーワードの割合が他に比べ少な過ぎたことが原因だと思われる。

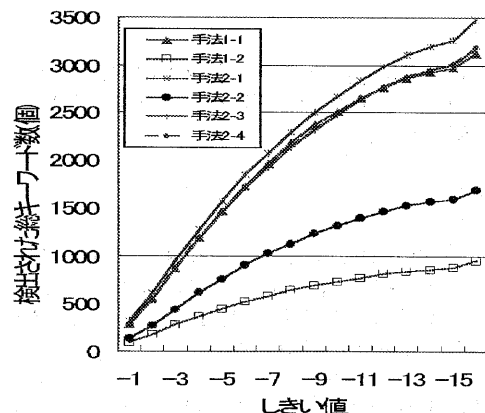


次に図1で示した手順で実際に話題決定実験を行った。各手法で最も決定率の高かったしきい値での結果を図4に示す。手法1-2を除いて、どの手法も同じような話題決定率となった。表2に手法1-1による話題決定の結果を示す。他の手法でもほぼ同様の傾向を示した。これからわかるように、話題認識結果がほとんど社会に偏ってしまい、もともと評価データの社会記事が占める割合が約69%であったため、どの手法も同じように70%弱の話題認識率を得たと考えられる。DF値を考慮することにより話題決定率の向上を期待していたが、図4のような結果となり、キーワード選択法の優劣をつけることができない結果となった。

次に表1に示すように、どの手法も出現キーワードの検出数が550から650単語(1話題あたり約6~7単語)までの間で、最も高い話題決定率を得ている。このことからキーワード数(今回は3000単語)が多すぎることが考えられる。キーワード数を減らせば、必然的に湧き出し誤り数も減少するので、より適正なキーワード数の考慮が必要である。参考までに各しきい値での出現単語の検出数を図5に示す。手法2-3と手法2-4は、ほぼ同じ結果となったので重なって表示されている。

表2 話題決定の結果(例:手法1-1)

	社	経	国	ス	科	芸	家	文
社	49	4	-	-	-	-	1	1
経	3	3	2	-	-	-	-	-
国	13	-	6	-	-	-	-	-
ス	2	-	-	-	-	-	-	-



次に手法1-1ではテキストデータで約80%の話題決定率が実際にスポットティングを行った実験で約65%と下がってしまい、逆に手法2-2では、約51%から69%へと上がってしまっている。テキストデータでの実験と音声データでの実験の結果が同じ程度の話題決定率であれば、湧き出し誤りの影響はないと考えられるが、これらのことより湧き出し誤りの影響が大きいことがわかる。

5 終わりに

本研究では、ワードスポットティングに基づいた話題決定を行うときに用いるキーワードセットの決定法を、ワードスポットティングにおける検出率/湧き出し誤り率および話題決定率の観点から評価を行った。結果として、

- ・ DF 値を考慮したキーワードセットを用いることによりスポットティング性能は向上した。
- ・ 話題決定では DF 値を考慮したキーワードセットの有効性は見られなかった。
- ・ 評価データに出現する単語のカバー率が低いと話題決定率も低くなった。

ことがわかった。このことより湧き出し誤りをいかに抑制するかということとワードスポットティング自体の性能向上が必要であると考えられる。

今後は、より良いキーワードセットの選出と、適正なキーワード数、湧き出し誤りの抑制の方法などを検討し、再実験を行う予定である。

謝辞

本研究では、キーワードセットの学習に毎日新聞 CD-ROM ('91年1月~'94年9月)を使用した。

参考文献

[1] 恒川 俊克、山下 洋一、溝口 理一郎 “キーワードスポットティングに基づく話題決定”、情処研報、SLP20-11、pp61-68 (1998-2)。

[2] 鷹尾 誠一、緒方 純、有木 康雄 “ニュース

音声の記事分類におけるキーワード選択法の比較”、情処研報、SLP22-15、pp75-82 (1998-7)。

[3] K. Ohtsuka, T. Matsuoka 他, “TOPIC EXTRACTION MULTIPLE TOPIC-WORDS IN BROADCAST-NEWS SPEECH”、ICASSP98、pp329-332 (1998)。

[4] 山下 洋一 “ワードスポットティングにおける湧き出し誤りの予測” 音語論集、1-1-11、pp21-22 (1998-9)。

[5] 山下 洋一, “音声認識シミュレーションに基づくワードスポットティング精度の予測”、信学技報 SP98-95、pp65-70 (1998-11)。