

ニュース音声に対するトピックセグメンテーションと分類

鷹尾 誠一 緒方 淳 有木 康雄

龍谷大学 理工学部

〒 520-2194 大津市瀬田大江町横谷1-5

Tel: 077-543-7427

E-mail: {tail,ogata}@arikilab.elec.ryukoku.ac.jp, ariki@rins.st.ryukoku.ac.jp

あらまし ビデオ・オン・デマンドを目指したニュースデータベースを構築するには、ニュース記事を話題(トピック)毎に分類する必要がある。しかし、従来の記事分類方法は、予め、記事内容毎に分割された記事を対象に行っていた。このため、ニュース記事のように複数の記事が連続して、その記事境界が未知である場合には、そのままでは適用することができなかった。そこで、本研究では、連続ニュース音声におけるトピック分類手法(トピックセグメンテーション)を提案する。

キーワード : 記事分類、教師ありトピックセグメンテーション、教師なしトピックセグメンテーション

キーワード

Topic Segmentation and Classification to News Speech

Seiichi Takao Jun Ogata Yasuo Arika

Faculty of Science and Technology, Ryukoku University

1-5 Yokotani, Oe-cho, Seta, Otsu-shi, 520-2194 Japan

Tel: +81-77-543-7427

E-mail: {tail,ogata}@arikilab.elec.ryukoku.ac.jp, ariki@rins.st.ryukoku.ac.jp

Abstract In order to construct a news database with a function of video on demand (VOD), it is required to classify a news articles into topics. So far, the article classification has been done to the articles, segmented in advance from the continuous news speech. However, in the case where the articles are not segmented, it was difficult to apply the classification to the continuous news speech. From this viewpoint, we propose a segmentation and classification method of the news articles from the continuous news speech.

Key words : article classification, supervised topic segmentation, unsupervised topic segmentation

key words

1 はじめに

近年、放送の多チャンネル化により、多くのニュース番組が放映されるようになった。これを受けて、視聴者には知りたいニュースだけを見たいという要求が生じている。この要求に対応するには、ニュース記事を分類して、データベースを構築しておく必要がある。しかし、人手でニュース記事を分類することは不可能であり、機械によるニュース記事の自動分類が望まれる。この点から、ニュース音声に対する話題分類、話題同定の研究が行われてきた。[1]-[6]。

しかし、これらの方法では、ニュースが予め記事毎に分割されていることを前提としている。連続したニュース番組を対象とする場合には、記事内容に基づいてニュース番組を記事毎に分割しなければならない。これはトピックセグメンテーションと呼ばれている。[7]-[9]。

トピックセグメンテーションでは、トピックの切目をニュース番組中で見つけるために、トピックに対してキーワードを学習しておく必要がある。これは、各記事に対して、トピックを教師データとして教えておく必要があり、教師ありトピックセグメンテーションと呼ばれている。これに対して、本研究では、トピックに対するキーワードの学習を必要としない教師なしトピックセグメンテーションを提案する。

この理由は、教師ありトピックセグメンテーションでは、分類されるトピックを教師データとして教えておく必要があり、セグメンテーションできるトピック数が固定されてしまうからである。例えば、政治・経済・国際のトピックを学習したとすると、セグメンテーションの際には、これらの3分類にしかセグメンテーションすることができない。更に、同じトピックの記事が連続して続いている場合は、本来は記事の切目であるにもかかわらず、教師ありトピックセグメンテーションでは、この切目を検出することができない。これに対して、教師なしトピックセグメンテーションでは、セグメンテーションの際に、トピック数を限定しないので、前述の教師ありトピックセグメンテーションの問題点を解決できるからである。

本稿では、まず最初に記事分類、次に教師ありトピックセグメンテーション、最後に教師なしトピックセグメンテーションについて述べる。

2 記事分類の概要

記事分類は図1に示すように、次の手順で行う。

1. RWCから提供されている形態素解析された93年の毎日新聞の記事10,000件を、朝日新聞92年度の方

類表索引[10]を用いて10トピックに分類する。分類方法は[2]に述べる従来法で行った。

2. 1の分類した結果を正解として、トピック毎にキーワードの選択、キーワードと分類トピックとの関連度の計算を行う。
3. ニュース音声記事をディクテーションした結果からキーワードを抽出する。次に、2で計算したキーワードの分類トピックに対する関連度を用いて、ニュース記事を分類する。

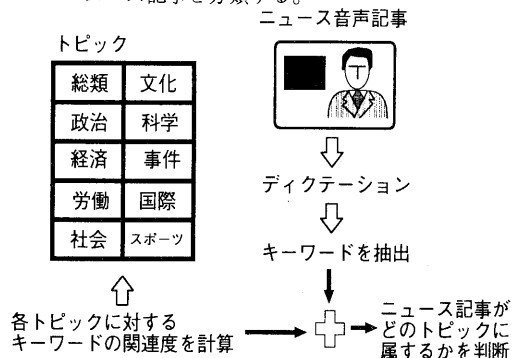


図 1: 記事分類の概要

3 新聞記事からのキーワード選択法

新聞記事中の単語 w_i と、その記事が属するトピック t_j との間で関連度 γ_{ij} を求め、閾値処理を行って、関連度 γ_{ij} の大きい単語をキーワードとして選択する。以下には、単語 w_i とトピック t_j との関連度 γ_{ij} を求める幾つかの方法を示す。

3.1 χ^2 値

χ^2 検定で用いられる χ^2 値は、トピック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_i の出現確率は全トピックを通じて等しいという仮説を設定し、この仮説に基づいて、単語 w_i のトピック t_j における予測頻度 m_{ij} を計算する。また、各単語 w_i について、トピック t_j における頻度 x_{ij} を求める。この x_{ij} と m_{ij} を基に、式(1)に従って χ^2 値を求める。もし、この χ^2 値が十分大きな値になれば、特定のトピックに偏って現れる単語と言う事になり、トピックの識別に有効な単語と見なせる。

$$\chi_{ij}^2 = \frac{(x_{ij} - m_{ij})^2}{m_{ij}} \quad (1)$$

$$m_{ij} = \frac{\sum_{j=1}^T x_{ij}}{\frac{W}{T}} \times \sum_{i=1}^W x_{ij}$$

- W : 異なり単語数
 T : トピック数
 x_{ij} : 単語 w_i のトピック t_j における頻度
 m_{ij} : 単語 w_i のトピック t_j における予測頻度

3.2 相互情報量

単語 w_i とトピック t_j との相互情報量 $i(t_j; w_i)$ とは、単語 w_i を知ることで、トピック t_j に関して得られる情報量のことであり、式(2)で表される。

$$\begin{aligned}
 i(t_j; w_i) &= i(t_j) - i(t_j|w_i) \\
 &= -\log P(t_j) + \log P(t_j|w_i) \\
 &= \log \frac{P(t_j, w_i)}{P(t_j)P(w_i)} \quad (2)
 \end{aligned}$$

また、相互情報量は文献[1]に示すように、 χ^2 値と同義であると言える。

3.3 TF-IDF

TF-IDFは式(3)で表され、単語 w_i が記事 a_k に現れる回数が高ければ高いほど、TF(Term Frequency)が高くなり、単語 w_i が現れる記事数が少なければ少ないほど、IDF(Inverse Document Frequency)が高くなる。したがって、TFは頻度の高い単語という性質を表し、IDFはそのトピックに偏って現れる単語という性質を表している。

$$\begin{aligned}
 TF \cdot IDF &= TF(w_i, a_k) \cdot IDF(w_i) \quad (3) \\
 TF(w_i, a_k) &= \text{単語 } w_i \text{ が記事 } a_k \text{ に現れる回数} \\
 IDF(w_i) &= \log \frac{\text{索引対象の全記事数}}{\text{単語 } w_i \text{ が現れる記事数}}
 \end{aligned}$$

本研究では、単語 w_i とトピック t_j との関連度を求め、これを閾値処理することでキーワードを選択している。このため、単語 w_i とトピック t_j との関連度を求める必要があるが、TF-IDFには本来、トピックという概念がない。そこで、式(4)のようにTF-IDFを変形してキーワードを選択した。

$$\begin{aligned}
 TF \cdot IDF &= TF(w_i, t_j) \cdot IDF(w_i) \quad (4) \\
 TF(w_i, t_j) &= \text{単語 } w_i \text{ がトピック } t_j \text{ に現れる回数} \\
 IDF(w_i) &= \log \frac{\text{全トピック数(10)}}{\text{単語 } w_i \text{ が現れるトピック数}}
 \end{aligned}$$

3.4 負の値をもつ χ^2 値

従来の χ^2 値では、単語 w_i がトピック t_j で発生する頻度 x_{ij} と予測頻度 m_{ij} の差を2乗して求めている。このため $x_{ij} < m_{ij}$ の場合であっても、 χ^2 値は正の値として

計算され、単語 w_i はトピック t_j において偏りがあると判断される。この問題を解決するために文献[3]と[6]では、式(5)に示すように改良方法が提案されている。

$$x_{ij}^2 = \frac{(x_{ij} - m_{ij}) \cdot |x_{ij} - m_{ij}|}{m_{ij}} \quad (5)$$

3.5 重み付き相互情報量

従来の相互情報量では、単語の発生頻度が小さくても、単語とトピックの依存度(共起関係)が高ければ大きな値を示してしまう。この問題を解決する方法として文献[6]では式(6)に示すように、従来の相互情報量に単語 w_i とトピック t_j の同時確率を重みとしてかける方法が提案されている。この値は、単語 w_i がトピック t_j に従属すると見た場合の確率分布と、独立すると見た場合の確率分布間のダイバージェンスとなっている。

$$I(w_i; t_j) = P(w_i, t_j) \cdot \log \frac{P(w_i, t_j)}{P(w_i)P(t_j)} \quad (6)$$

3.6 相対相互情報量

式(2)に示す相互情報量は、 $P(w_i|t_j) < P(w_i)$ の場合に負の値となる場合がある。また、相互情報量にはトピック間の相対的な関係が考慮されていない。そこで式(7)に示すように、 $P(w_i)$ の代わりに、全てのトピックの中で単語 w_i が出現する確率が最小のもの、すなわち、 $\min_k P(w_i|t_k)$ を用いる。[1]。こうすることで、単語 w_i のトピック t_j における出現確率のダイナミックレンジを大きくすることができる。この値は、相互情報量、あるいは χ^2 値における基本表現 x_{ij}/m_{ij} の分母を、 $\min_k P(w_i|t_k)$ として強調したものになっている。

$$\frac{p(w_i|t_j)}{\min_k p(w_i|t_k)} \quad (7)$$

4 記事の分類方法

1つの記事 x が与えられると、記事中のキーワード w_i を全て抽出する。このキーワード w_i とトピック t_j との関連度 γ_{ij} を基に、キーワード w_i がトピック t_j の分類に寄与する割合 C_{ij} を計算する。最後に、1つの記事中に含まれているキーワードの発生回数 x_i を、式9のように正規して N_i を求める。最後に、分類寄与率 C_{ij} を要素に持つベクトル v_j と、 N_i を要素に持つベクトル x の内積を記事 x とトピック t_j との類似度として求める。この類似度が一番大きいトピックに記事 x を分類する。

$$S(x, t_j) = x \cdot v_j = \sum_i N_i \cdot C_{ij} \quad (8)$$

以下には、キーワード w_i のトピック t_j に対する分類寄与率 C_{ij} について、幾つかの計算方法を示す。

4.1 関連度に基づく類似度

キーワード w_i とトピック t_j との関連度 γ_{ij} を式 (9) のように正規化して分類寄与率 C_{ij} を求め、記事を分類する方法である。この分類方法を C^1 とする。

$$C_{ij} = \frac{\gamma_{ij}}{\sum_j \gamma_{ij}} \quad (9)$$

この方法は、キーワード w_i のトピック t_j に対する関連度を、キーワード間で比較可能にする効果がある。これは文献 [1] において提案された手法である。

4.2 キーワードの出現回数に基づく分類寄与率

キーワード w_i がトピック t_j において出現した回数 x_{ij} を基に、分類寄与率 C_{ij} を式 (10) のようにして求め、記事を分類する方法である。この分類方法を C^2 とする。

$$C_{ij} = \frac{N_{ij}}{\sqrt{\sum_i N_{ij}^2}} \quad (10)$$

この分類方法は、式 8 より、

$$\begin{aligned} S(x, t_j) &= x \cdot v_j = \sum_i N_i \cdot C_{ij} \\ &= \|x\| \|v_j\| \cos \theta_j = \|x\| \cos \theta_j \end{aligned}$$

であることから、単純類似度法と同じである。

5 ニュース音声ディクテーション結果に対する記事分類実験

5.1 ディクテーション

5.1.1 実験条件

用いた言語モデルは、毎日新聞 CD-ROM 版の 45ヶ月分 (91年1月～94年9月) の記事から学習したものである。語彙数 20K の back-off bigram で、back-off smoothing には witten-bell の推定を用いている。bigram に対する cut-off は 1 とした。

音響モデルは、男性不特定話者 HMM で、単語間の音素文脈依存も考慮した cross-word triphone モデルである。学習には、日本音響学会新聞記事読み上げ音声コーパスのうち、男性話者 137 人分の 21782 発話を用いた。音響特徴量には 39 次元の特徴パラメータ (12 次元の MFCC とパワー、およびそれぞれの Δ 、 $\Delta\Delta$ 係数) を用いた。

5.1.2 連続音声認識実験

評価用音声データには、93年、94年のNHK1時のニュース55記事分(総計1.13時間、1記事平均71秒)を用いた。bigram 言語モデルの学習データとは時期的に closed なデータである。55記事のデータを自動的に1発話毎に区切った総発話数409のデータに対してディクテーションを行った。無音区間が1秒以上連続した場合を発話の切目として、ニュースを発話単位に切り出している。

記事分類用音声データの諸元を表1に、ディクテーション結果を表2に示す。表2のキーワード正解率と適合率は次のように定式される。ただし、ディクテーションの結果得られた真のキーワード数を x 、ニュース音声記事を手書きおこしたテキスト中の全キーワード数を y 、ディクテーションの結果得られた全キーワード数を z としている。

$$\text{キーワード正解率} = x/y \quad (11)$$

$$\text{キーワード適合率} = x/z \quad (12)$$

表 1: 記事分類用音声データ

話者数	3名
総発話数	409
perp	78.3
20K 未知語率	0.8%

perp: test-set perplexity

表 2: ディクテーション結果

単語正解率	単語正解精度	単語誤り率
85.6%	80.3%	19.7%
キーワード正解率	キーワード正解精度	キーワード誤り率
92.1%	81.7%	18.3%

5.2 記事分類実験結果

実験は、6つのキーワード選択法と2つの分類方法の組み合わせ12種類を行った。以下では、各キーワード選択法について一番結果が良かったときの結果だけを表3に示す。表3の Chi2 は χ^2 値、Mutual は相互情報量、TF-IDF は TF-IDF、Chi2-Negative は負の値を持つ χ^2 値、Mutual-W は重み付き相互情報量、Mutual-R は相対相互情報量である。

表 3: 記事分類実験結果

キーワード選択法	分類方法	キーワード数	記事分類率
Chi2	C^1	6294	81.8%
Mutual	C^1	13755	81.8%
TF-IDF	C^2	366	78.1%
Chi2-Negative	C^1	13755	81.8%
Mutual-W	C^1	13755	81.8%
Mutual-R	C^1	24940	83.6%

表3に示すように、 χ^2 値、相互情報量、負の値を持つ χ^2 値、重み付き相互情報量、相対相互情報量は C^1 と相性がよく、TF-IDF は C^2 と相性が良いことがわかる。なお、詳細な実験結果は文献 [1] に示されている。

6 教師ありトピックセグメンテーションの概要

教師ありトピックセグメンテーションは、図2に示すように次の手順で行う。

1. キーワードとトピックとの関連度を、3 節で述べた方法によって予め求めておく
2. ニュース音声をディクテーションしてキーワードを抽出する。
3. 分析区間内でキーワードの頻度分布を求める。
4. キーワードの頻度分布をベクトル x 、トピック t_j に対するキーワードの関連度を要素に持つベクトルを v_j とすると、4節の式 (8) により、分析区間とトピックの類似度を求める。
5. 分析区間をずらしながら、3~4を繰り返すことによって、各トピックの時間関数(トピック関数)を求め、トピック関数が最大となる区間をトピック区間(記事)として切り出す。

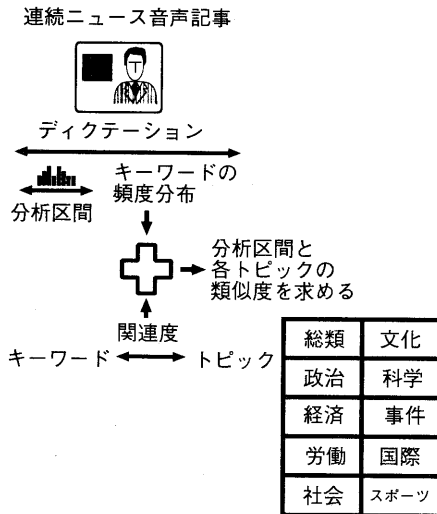


図2: 教師ありトピックセグメンテーションの概要

7 教師ありトピックセグメンテーションの実験

5.1節で述べた方法により、連続ニュース音声をディクテーションしておく。この結果に対して、教師ありトピックセグメンテーションを行い、記事に分割する実験を行った。

7.1 教師ありトピックセグメンテーションの実験結果

実験は、6つのキーワード選択法に対して、記事分類で相性の良かった分類方法を用いて6種類のセグメンテーション実験を行った。また、分析区間長は、1発話から10発話まで変化させて実験を行った。表4に、各キーワード選択法について、一番良かったときの結果を示す。図3は、分析区間長を1発話としたときに、キーワードの数を変えた場合の区間正解率の変化を示している。

表4: 教師ありトピックセグメンテーションの実験結果

キーワード選択法	分析区間長	キーワード数	区間正解率
Chi2	1	24940	61.81%
Mutual	1	20923	63.63%
TF-IDF	1	745	54.54%
Chi2-Negative	1	24933	67.27%
Mutual-W	1	13755	60.00%
Mutual-R	1	24940	50.90%

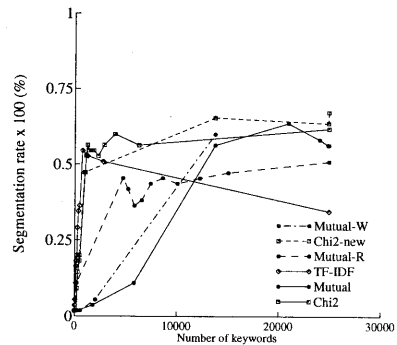


図3: 分析区間長が1発話のときの区間正解率のグラフ

区間正解率とは、切り出された区間の始めと終りを、正しくセグメンテーションできた区間の割合である。例えば、次の図4では、政治の終りと経済の始めが切り誤っており、トピック区間としては国際のみが正しく切り出されている。従って、区間正解率は $\frac{1}{3}$ となる。

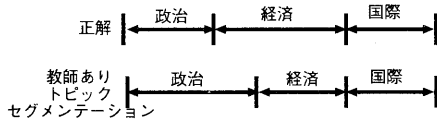


図4: 教師ありトピックセグメンテーションの評価方法

表4に示されるように、教師ありトピックセグメンテーションにおいては、負の値を持つ χ^2 値が良いことがわかる。これは、負の値を持つ χ^2 値はキーワードの関連度のレンジが広く、負の値も持つことによって、従来の χ^2 値の弱点を克服しているためだと考えられる。

また、表3の記事分類に比べて、キーワード数が多くなければ十分な区間の正解率が得られないことがわかる。分析区間が記事単位から発話単位となることで、分析区間長が短くなったため、キーワード数が多く必要になったものと考えられる。

8 教師なしトピックセグメンテーションの概要

8.1 教師なしトピックセグメンテーションの手順

教師なしトピックセグメンテーションは図5に示すように次の手順で行う。

1. ニュース音声をディクテーションする。
2. ディクテーション結果において分析区間をきめ、全ての分析区間において、単語(名詞)の頻度分布を求める。
3. 分析区間毎に、単語(名詞)の頻度分布から単語(名詞)の重要度を決定する。
4. 単語(名詞)の重要度を基に、分析区間毎にトピックベクトルを作成する。
5. 隣接する分析区間対において、トピックベクトルの類似度を求め、類似していればトピックが継続していると判断し、類似していなければトピックの境界と判断する。

連続ニュース音声記事

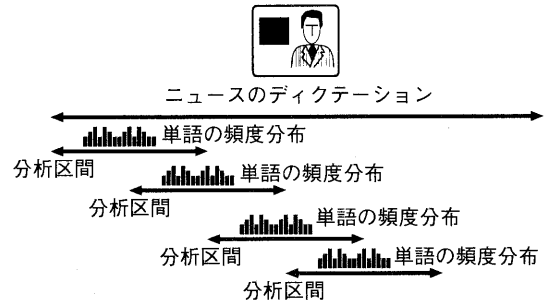


図5: 教師なしトピックセグメンテーションの概要

8.2 分析区間における単語の重要度

1つの分析区間において、単語の重要度を決定するには、単語 w_i が全ての分析区間で均等に出現すると仮定した場合に対して、どの程度偏りがあるかを求めて判定すればよい。従って、3.1節で述べた χ^2 値を用いて、分析区間における単語の重要度を決定することができる。ただし、 x_{ij} と m_{ij} は、分析区間 t_j における単語 w_i と予測頻度に変更している。

$$\chi_{ij}^2 = \frac{(x_{ij} - m_{ij})^2}{m_{ij}} \quad (13)$$

$$m_{ij} = \frac{\sum_{j=1}^T x_{ij}}{W \sum_{j=1}^T x_{ij}} \times \sum_{i=1}^W x_{ij}$$

W : 異なり単語数

T : 分析区間数

x_{ij} : 単語 w_i の分析区間 t_j における頻度

m_{ij} : 単語 w_i の分析区間 t_j における予測頻度

8.3 分析区間のトピックベクトル

分析区間毎に単語の重要度を閾値処理して、重要度の高い単語(名詞)だけを取り出し、トピックベクトルを作成する。トピックベクトルの作成方法として、次の2つの方法がある。

- 単語の頻度を成分とするベクトルを作成する。
- 単語の重要度を成分とするベクトルを作成する。

8.4 トピックベクトルの比較

分析区間 i におけるトピックベクトルを

$$X_i = (x_{1i}, x_{2i}, \dots, x_{ki}, \dots, x_{ni})$$

分析区間 j におけるトピックベクトルを

$$X_j = (x_{1j}, x_{2j}, \dots, x_{kj}, \dots, x_{nj})$$

とすると、その類似度は、

$$\cos \theta = \frac{(X_i, X_j)}{\|X_i\| \|X_j\|} \quad (14)$$

で表される。 $\cos \theta$ が 1 に近ければ近いほど、類似度が高い。

9 教師なしトピックセグメンテーションの実験

5.1節で述べた方法により、連続ニュース音声をディクテーションしておく。この結果に対して、教師なしトピックセグメンテーションを行い、記事に分割する実験を行った。

9.1 教師なしトピックセグメンテーションの実験結果

実験は、分析区間長を 1 発話から 10 発話まで変化させて行った。一番適合率が良かったときの結果を表 5 に、類似度の閾値によって正解率と適合率が変化する様子を図 6~7 に示す。

なお、表と図中の Frequency は単語の頻度を成分とするトピックベクトル作成方法、Keyword-Value は単語の重要度を成分とするトピックベクトル作成方法である。

図 6~7 からわかるように、教師なしトピックセグメンテーションでは境界の適合率がかなり低く、50%弱が最高である。適合率が低い原因としては、境界の付近では複数のトピックが混在するために、トピック決定が不安定になるので、境界付近で複数のセグメントを生成してしまうためである。

表 5: 教師なしトピックセグメンテーションの実験結果

単語の重要度の決定方法	χ^2 値	χ^2 値
トピックベクトルの作成方法	Frequency	Keyword-Value
分析区間長	3 発話	5 発話
単語の重要度の閾値	10	0
類似度に対する閾値	0.02	0.01
境界正解率	85.2%	90.7%
境界適合率	46.5%	40.2%

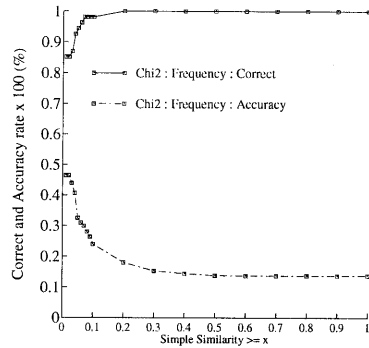


図 6: 分析区間長が 3 発話、単語の重要度が 10 のときの境界正解率と適合率のグラフ

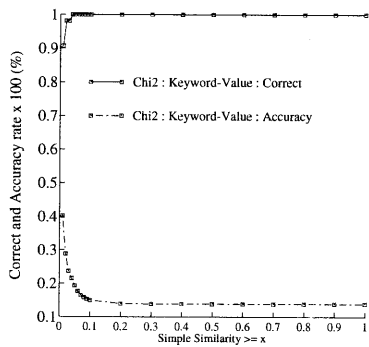


図 7: 分析区間長が 5 発話、単語の重要度が 0 のときの境界正解率と適合率のグラフ

10 終りに

本研究では、NHK ニュース音声のディクテーションに対して、記事分類、教師ありトピックセグメンテーション、教師なしトピックセグメンテーションの実験を行った。記事分類では各キーワード選択法に優劣がつけがたかったが、教師ありトピックセグメンテーションにおいては、負の値を持つ χ^2 値が優位性を示した。これは、先にも述べたように、キーワードの関連度のレンジが広く、かつ、従来の χ^2 値の弱点を克服しているためだと考えられる。

また、今回行った教師なしトピックセグメンテーションの実験は、単語の重要度の決定方法を χ^2 値でしか行っていない。今後は相互情報量、TF-IDF、負の値を持つ χ^2 値、重み付き相互情報量、相対相互情報量で教師なしトピックセグメンテーションの実験をし、これらの手法が教師なしトピックセグメンテーションにおいて、どのような優劣を示すかを比較・検討したい。

更に、教師ありトピックセグメンテーションと教師なしトピックセグメンテーションを統合することによって、

トピックセグメンテーションの精度向上を計りたい。

参考文献

- [1] 鷹尾誠一, 緒方淳, 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”, 音声言語情報処理, **SLP98-22**, pp.75-82 (1998-07).
- [2] 櫻井光康, 有木康雄: “キーワードスポッティングによるニュース音声の分類と索引付け”, 信学技法, **SP96-66**, pp.37-44 (1996-11).
- [3] 藤井洋一, 今村誠, 高山泰博, 鈴木克志: “共起情報を利用した新聞記事の自動分類結果の分析・評価”, 情報処理学会第55回(平成9年後期)全国大会, pp.(3-212)-(3-213).
- [4] 高木幸一, 桜井直之, 岩崎淳, 古井貞熙: “ニュース音声を対象とした言語モデルと話題抽出の検討”, 信学技報, **SP98-33**, pp.73-80 (1998-06).
- [5] 岩崎淳, 古井貞熙: “ニュース音声からの話題抽出法の検討”, 日本音響学会 平成10年度秋季研究発表会, pp.27-28.
- [6] K. Ohtsuki, T. Matsuoka, S.Matsunaga, S.Furui: “TOPIC EXTRACTION MULTIPLE TOPIC-WORDS IN BROADCAST-NEWS SPEECH”, ICAS-SP98, pp.329-332 (1998).
- [7] 鷹尾誠一, 緒方淳, 有木康雄: “ニュース音声に対するトピックセグメンテーションの検討”, 日本音響学会 平成10年度秋季研究発表会, pp.157-158.
- [8] 鈴木良弥, 福本文代, 関口芳廣: “ニュース文を対象とした話題毎のセグメンテーション”, 日本音響学会 平成8年度春季研究発表会, pp.185-186.
- [9] 恒川俊克, 山下洋一, 溝口理一郎: “キーワードスポッティングに基づくニュース音声の話題分類”, 音声言語情報処理, pp.61-68 (1998.2.6).
- [10] “朝日新聞記事データベース分類表索引”, 朝日新聞社ニューメディア本部, (1992).