

テキスト構造を利用した主題の推定について

野本 忠司

(株)日立製作所 基礎研究所

〒 350-03 埼玉県比企郡鳩山町赤沼 2520

email: nomoto@harl.hitachi.co.jp

松本 裕治

奈良先端科学技術大学院大学

〒 630-01 奈良県生駒市高山 8916-5

email: matsu@is.aist-nara.ac.jp

あらまし

ネットワーク上に流通する電子出版物はインターネットの普及にともない、今後さらに増加し、情報の効率的な選別の必要性がますます高まってくると予想される。このような背景のもと、本稿ではテキストの主題(トピック)を推定する方法について述べる。本手法は、基本的には従来の文書分類の技術に基づいたものであるが、本稿では特にテキストの構造を利用した主題の推定方式について実際のコーパスデータをもとに検討した。CD-ROM版日本経済新聞(1992年1~6月)の約4万件の記事を利用し実験をおこなった結果、最も良い場合でテキストの構造的な特徴を利用しない方式に比べて、平均で20%程度の精度の向上が確認された。

キーワード 日本語、ディスコース、テキスト分類、主題抽出、テキストの構造化

Exploiting Text Structure for Topic Identification

Tadashi Nomoto

Advanced Research Laboratory, Hitachi Ltd.

2520 Hatoyama Saitama 350-03 Japan

Yuji Matsumoto

Nara Institute of Science and Technology

8916-5 Takayama Ikoma Nara, 630-01 Japan

Abstract

The paper demonstrates how information on text structure can be used to improve the performance on the identification of topical words in texts, which is based on a probabilistic model of text categorization. We use texts which are not explicitly structured. A text structure is identified by measuring the similarity between segments comprising the text and its title. It is shown that a text structure thus identified gives a good clue to finding out parts of the text most relevant to its content. The significance of exploiting information on the structure for topic identification is demonstrated by a set of experiments conducted on the 19Mb of Japanese newspaper articles. The paper also brings concepts from the rhetorical structure theory (RST) to the statistical analysis of a text structure. Finally, it is shown that information on text structure is more effective for large documents than for small documents.

key words **Japanese, discourse, text categorization, topic identification, text structuring**

1 はじめに

元来、人間はスキミング(飛ばし読み)、スキヤニング(探し読み)などの読書能力を駆使して書物の取捨選択をおこなってきたが、この能力をコンピューターでより効果的に支援することは、大量の情報が電子化され氾濫している現代にあっては、ますます必要になってくると考えられる。本稿では、情報の選択また一般的な読書の支援機能の一つとして、主題(トピック)の抽出方法について検討する。

テキストの主題とは、その内容をよく表わしている表現であり、一般的には名詞であることが多い。主題に関しては、言語学でも定義がいまだに確立されていないが、本稿では便宜的にタイトル(見出し)に表われている名詞を主題と呼ぶことにする。主題の抽出は技術的には情報検索におけるインデクシング(単語の重みづけ)に近く、本研究も基本的にはテキスト分類という情報検索の一分野で提案された手法に基づいている。ただし、一般にテキスト分類が文書を予め決められたカテゴリーに割り振ることを目的としているのに対し、本方式では、テキストをそれ自身に表れている単語(名詞)をカテゴリーとして分類する。情報検索では通常、文書のサイズが大きくなると精度が低下することが知られている。大文書には本筋とは関係ないさまざまな小トピックが現れ、これらがノイズとなって検索を混乱させるからである。このような問題に対し、検索、特に Passage Retrieval という分野で文、段落、節、章など書式を含めたテキストの構造を利用した検索手法が研究されている。Wilkinson (1994) は書式構造以外に *abstract*, *purpose*, *supplementary* など意味的な概念の導入を提案しており、また Salton *et al.* (1993) では節、段落を一文書とした検索実験を行い、全文検索の場合より良い結果が得られたと報告している。さらに Hearst (1994) では統計的な手法により文書をタイルというまとまりに分割し、タイルを利用した部分検索を提案し、全文検索より優位であることを示した。また、Callan (1994) では部分テキスト (passage) と全文を共に使った折衷方式を提案している。以下、第2節では主題推定の技術的な詳細を述べる。さらに当該方式の問題点(第3節)、解決方法の提案(第4節)を経て、第5節と第6節で実験の方法と結果を報告する。

2 主題推定

本稿における主題推定とはタイトルなどの情報を利用せず本文の主題を推定することである。前節でも述べたように本稿ではテキストの見出しに表われる名詞主題を主題と定義する。大きな理由は、評価を簡便にするためである。したがって、主題を推定するという作業は全く見出しを見ずにテキストの見出し語を本文から当てることである。無論、主題として見出し語だけではなくその類義語も正解にすることもできるが、今回は厳密に見出し語のみに限った。

本稿で導入する手法は基本的にはテキスト分類 (text categorization) のそれである。テキスト分類とは、情報検索の一つの研究分野であるが、その目的はテキスト(文書)を予め用意した項目 (category) に正しく分類するアルゴリズムの開発である。通常、テキスト分類は以下のステップで構成される。

Step 1 用意された分類項目それぞれについて、文書がその項目に分類される確率あるいは確率に相当する得点を計算する。

Step 2 次に得点を基に(適当な閾値に基づいて)、実際に文書に項目を付与する。

技術的な細部に入る前に、具体例で若干説明しておく。まず、図1を参照されたい。図は実際訓練データに用いた記事であり、見出しと記事本文より成っている。また、記事は形態素解析プログラム JUMAN を用いて解析した。括弧内は JUMAN (Matsumoto *et al.*, 1994) が名詞と認定した形態素を表わす。¹ (以後、これらの形態素を「名詞」と呼ぶ。) ここで、主題候補語を potential topic、実際見出しに現われた名詞を actual topic と呼ぶ。図1の記事を d とすると potential topic $S(d)$ また actual topic $T(d)$ は 以下ようになる。また、 $W(d)$ は本文中に現われた名詞の集合である。ただし、 $S(d) \subseteq W(d)$ とするが、ここでは $W(d) = S(d)$ を仮定する。すなわち、本文に現われた名詞をすべて主題の候補語と考える。本稿後半で、 $S(d)$ の規定方法と主題推定の精度との関係について議論する。

テキスト分類では、予め用意した項目と文書との関係に注目するが、主題推定ではテキスト本文に表われた名詞と本文の関係に注目する。本稿では、以下の式で名詞の主題らしさを定義する。 c は主題候補語 (potential topic)、 d はテキストを表わす。 d はテキストの全

¹星印は解析エラーを表わす。

(仏)(銀)が(キエフ)に(駐在)(員)(事務所)

(仏)(銀)(大手)の(ソシエテ)(ジェネラル)は15日、(ウ*)(クラ*)(イナ*)の(首都)(キエフ)に(駐在)(員)(事務所)を(開設)すると(発表)した。すでに(キエフ)(市)(当局)の(許可)を得たと言う。

図 1: 名詞の抽出 (日経記事より)。上が見出し、下が本文。

$$T(d) = \{ \text{仏, 銀, キエフ, 駐在, 員, 事務所} \}$$

$$S(d) = \{ \text{仏, 銀, 大手, ソシエテ, ジェネラル, 15日, ウ-, クラ-, イナ, 首都, キエフ, 駐在, 員, 事務所, 開設, 発表, 市, 当局, 許可} \}$$

$$W(d) = \{ \text{仏, 銀, 大手, ソシエテ, ジェネラル, 15日, ウ-, クラ-, イナ, 首都, キエフ, 駐在, 員, 事務所, 開設, 発表, 市, 当局, 許可} \}$$

体あるは一部分であってもよいが、見出し情報を含んではならない。

$$\mathcal{L}(c|d) = \sum_{t \in S(d)} P(c|t)P(t|d)$$

c と d の関連性を直接観測することができないので、 t という仲介表現 (mediating term) を導入する。 t は単語、 n -gram、文節等、何であってもよい。ちなみに、本稿では主題候補語 (potential topic) を仲介表現として使った。さて、上式をベイズの定理を使って変形すると以下ようになる。

$$\mathcal{L}(c|d) = P(c) \sum_{t \in S(d)} \frac{P(t|c)P(t|d)}{P(t)}$$

それぞれの確率を以下で推定する。

$$\begin{aligned} P(c) &= N(D_c)/N(D) \\ P(t|c) &= f(t, D_c)/|D_c| \\ P(t|d) &= f(t, d)/|d| \\ P(t) &= f(t, D)/|D| \end{aligned}$$

ここで、 D は訓練用テキスト(記事)の集合。 D_c は表現 c を見出しに含む記事の集合。 $N(D)$ は D の総記事数。同様に、 $N(D_c)$ は D_c の総記事数。 $f(t, D_c)$ は D_c における表現 t の(本文中の)総頻度。 $f(t, d)$ は d における t の頻度。 $f(t, D)$ は D における t の総頻度。 $|D_c|$ は D_c の出現単語総数。 $|d|, |D|$ も同様。

3 問題

主題推定方式の現状での問題は記事が長くなると推定の精度が低下するという点である。これは $S(d) = W(d)$ とする仮定に起因すると考えられる。当然であるが、この仮定のもとでは、記事が長くなればそれだけ主題候補の数も増える。図2は実際に見出しに表われた名詞が本文中にどのくらいの割合を占めるかを示したグラフである。見出しの長さはテキストの長さに影響されず一定しているので、本文が長くなれば、見出し語はそれだけ無作為に抽出される可能性が低くなる。例えば、100では見出し語は本文中に占める割合は13%であるが、900になるとその割合が3%まで低下する。以下ではこの問題の解決法について検討する。

4 テキスト構造を利用した主題推定法

本節では、テキスト構造の情報を利用して前節の問題の解決を試みる。基本的な目標としては、テキスト構造を利用して重要個所とそうでない個所を識別し、potential topic $S(d)$ の圧縮することを考える。なお、今回データとして用いる新聞記事にはテキストの構造を示す明確な言語的な標識がないため、統計的な手法を使って構造を検出することにした。テキスト構造検出の統計的手法としてよく知られているもので、テキス

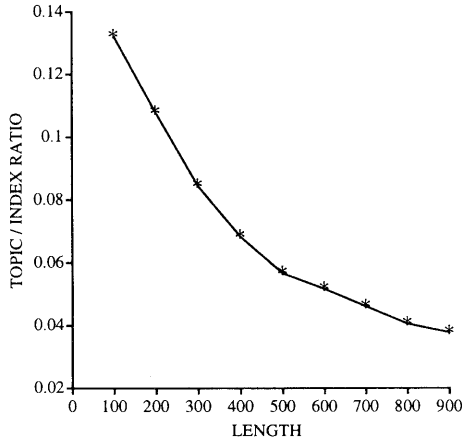


図 2: $T(d)$ と $W(d)$ の割合 (縦軸)。横軸は記事のサイズ (文字数)。データはテストコーパスの記事を利用した。100 は長さ 101 文字以上 200 文字以下の記事。200 は長さ 201 文字以上 300 文字以下の記事。他同様。

トをいくつかのブロックに分割して隣り合ったブロックの類似度をもとに構造を決めるという Hearst (1994) の方法があるが、この方式はテキストの重要部分の検出という目的には有効ではない。今回われわれが考えるのは、見出しとブロックの類似度によるテキストの重要部分の同定である。これは、見出しというのが記事のテーマであり、見出しに近いブロックがその記事の中心的な箇所であるという考えに立っている。

では、類似度の求めかたについて見ていく。本研究では以下の式で類似度を定義する。(Wilkinson, 1994) ここで、 h は見出し、また d は記事本文を表わすが、実際にはそれぞれ自身に現われる名詞の集合として考える。

$$SIM(h, d) = \sum_{i=0}^N ntf_{td} \cdot idf_i$$

N は h に現われる総語数、 ntf_{td} は d における t の正規頻度。定義は以下の通り。

$$ntf_{td} = \frac{tf_{td}}{\max_t tf_d}$$

tf_{td} は d における t の頻度、 $\max_t tf_d$ は d のなかの最頻度語の頻度。

$$idf_i = \frac{\log \frac{df}{idf}}{\log df}$$

tdf は d において t が現われるブロックの数。 df は d のなかのブロックの総数。実験では、見出しから単語

の重複を取り除き、ブロックの長さは 10 単語に固定した。 idf (inverse document frequency) は情報検索で用いられる文書単位の単語の重み付けのひとつであるが、特徴語ほど局所的な分布を持つという直観に基づいている。例えば、「桜」が 100 文書、「梅」が 10 文書に出現したとすると idf の値は「梅」のほうが高くなる。

この類似度の尺度を使って、訓練用記事について解析した結果が以下の通りである。(図 3) 手続きとしては、訓練コーパスを記事の長さ別にいくつかの組に分け、それぞれの組についてそこに含まれる記事を 10 単語単位のブロックに分割する。ちなみに 10 単語に満たないブロックも一応計算に対象に入れることにした。次に、その各ブロックに対して上記の類似度関数を使って、見出しとの類似度を計算した。

データは 1992 年 1 月から 1992 年 4 月まで日本経済新聞の記事を利用した。(日本経済新聞社, 1992) それぞれのグラフは記事の長さ別になっている。‘100-200’ は 100 から 200 文字長の記事を示す。他同様。また、‘t’ は関連サイズの記事数を示す。また、本稿後半で紹介する主題推定の実験では 100 から 1000 文字長の記事を対象にするので、ここでは 100 文字以下、1001 文字以上の記事については解析から除外している。グラフを見えることは、新聞記事では頭の部分に見出しとの最も類似しているブロックが集中しているということである。末尾では逆に極めて少くなっている。この結果から、われわれは、記事の末尾はあまり主題と関連性がなく、主題推定では無視することができるのではないかと考えた。次節ではこの考えを実験的に検証していく。

さて、図 3 のグラフの定性的な意味付けについて若干触れたい。この高い類似度の前方への偏りは、修辞理論的に説明するならば、新聞記事においては核 (nucleus) 要素が先頭に位置し、それに付加要素 (adjunct) が続くという構造になっているのではないかと考えられる。(Fox, 1987) 実際、例にあたってみる。例えば、図 1 の記事は、以下の 3 つの文から成っている。

- (1) 仏銀がキエフに駐在員事務所
- (2) 仏銀大手のソシエテジェネラルは 15 日、ウクライナの首都キエフに駐在員事務所を解説すると発表した。
- (3) すでにキエフ市当局の許可を得たと言う。

(1) は見出しであるが、残りの 2 文については、(2) が見出しとの関連から記事全体の中心的な主張、(3) が前

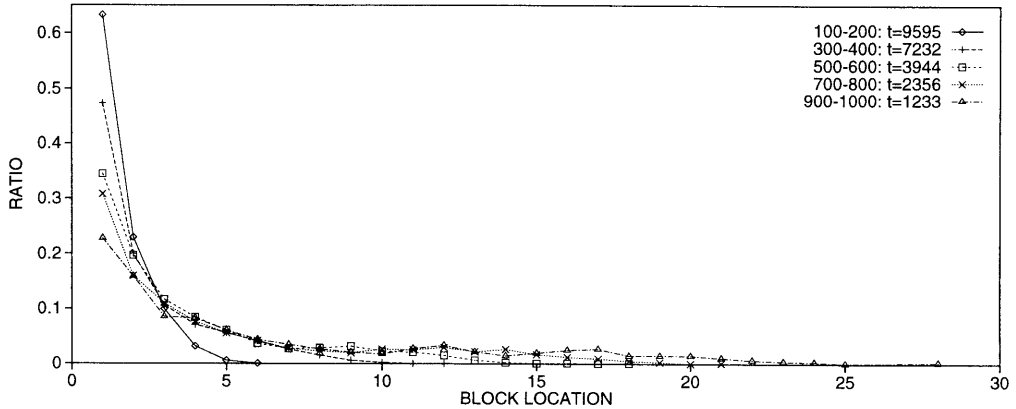


図 3: 類似度を利用したテキスト構造の解析。グラフの横軸はブロックの表われた記事の頭からの位置。縦軸はそれぞれのサイズにおいて各位置のブロックの類似度が最も高い場合の割合を示す。また、位置は最初のブロックを 0 番目として数える。例えば、5 ブロック目 ($x = 5$) は 40 番目の単語から始まり 49 番目の単語で終る。 t は各セットの記事数。

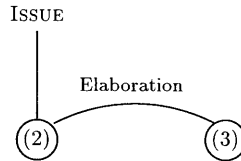


図 4: 図 1 の記事に対応する修辭構造

文に関する補足的な情報という解釈ができる。これを修辭構造の観点から表現すると、図 4 のようになる。記事本文は全体として ISSUE 構造をとり、核要素と付加要素が ELABORATION という関係に立つ。(Fox, 1987) また、他のいくつかの記事と先の類似度グラフの結果を参考にして考えると、多少の違いはあれ、一般的に記事のテキスト構造は図 5 のようになっているのではないかと推察される。

5 手続き

以下では、主題推定に全文を用いた場合(全文方式)と部分(部分方式)を用いた場合で、精度に差が見られるかどうか実験してみた。実験にはデータとして 43,253 のフルテキストの日経の記事(日本経済新聞社, 1992)を使った。記事は 92 年上半年(92 年 1 月～6 月)に掲載されたもので、1 月から 5 月までの 40,553 記事を訓練データ、6 月の 2,700 記事をテストデータ用にした。

訓練セット、テストセットともにもとのコーパスか

ら JUMAN を使って、名詞のみを抜き出して作った。次にテストセットを長さ別に九つのグループに分けた。各グループ 300 記事で構成されている。サイズは文字数で表わしている。テストセット 1 はサイズが 101 文字以上 200 文字以下の記事、テストセット 2 は 201 文字以上 300 文字以下という具合になっている。

さらに、部分方式では記事の先頭部分以外は切り捨てることにした。「先頭」部分をどう決めるかについては試行錯誤的に 3 つのやり方を試してみた。一つは、記事の長さに関りなく、先頭から一定の長さのブロックを切り出すやり方、二つ目は、記事の長さに比例して先頭部を切り取る方法、そして三つ目は、第一段落を先頭部分とする方法である。以下では、それぞれ、fixed-length model (FLM) proportional-length model (PLM) first-paragraph model (FPM) と呼ぶ。

6 実験と評価

以下の表 2、表 3 ではそれぞれ FLM、PLM 方式の結果を示してある。値は 'break-even point' つまり適合率 (precision) と再現率 (recall) が一致する最大値で

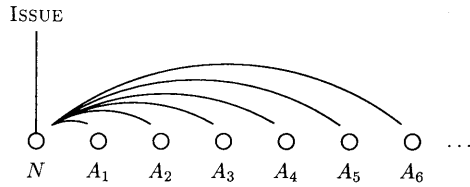


図 5: ニュース記事の修辞構造(予想)。「N」は核要素、「A」は補足要素。アークは要素間の関係を示す。

示している。ちなみに、適合率、再現率は以下で定義する。

$$\text{適合率} = \frac{\text{正解単語数}}{\text{見出し語として推定した単語数}}$$

$$\text{再現率} = \frac{\text{正解単語数}}{\text{実際の見出し語の数}}$$

ここで、**正解単語数**とは主題推定モデルが見出し語として推定した単語で実際の見出し語と一致した単語の数を指す。第2節で与えた推定モデルは「見出し語らしさ」について、主題候補単語に順位を付けるのが目的である。主題推定の作業は、さらに何点まで主題語として合格にするかということを決めなくては行けない。いわゆる、「thresholding」である。(Lewis, 1992) 一般的に、合格ラインを上げると、適合率は高くなるが再現率は落ちる。また、下げると、逆に適合率が落ち、再現率が高くなる。

表2と表3に戻る。表中の「+/-」は全文方式と比較して何%向上/劣化しているかを示している。また、*印の箇所は break-even point が存在しないので、代わりに最大再現率(括弧内)での適合率を示してある。固定方式(FLM)について、テキストの長さがブロック長に満たない場合はどこも切り捨てず、全体をとることにした。

結果をみると、FLM, PLM 方式いづれも全文方式に比べ、精度が向上しているのがわかる。特に、テキストのサイズが大きいのもの ($l > 200$) が小さいもの ($100 < l \leq 200$) に比べて、精度がよい。しかし、「切捨て」方式はテキストが長くなるにつれ伸びが鈍ってくる。例えば、表2で、 $i = 20$ の場合、つまり、記事の先頭20単語のブロックで主題推定のおこなうと、600-700の記事で精度の向上率はピークに達して、900-1000になるとピークに対して10%の精度低下が見られる。原因として考えられるのは、(図3を眺めてみると、) テキストのサイズが大きくなるにつれ類似度分布の偏向が緩やかになるという点である。つまり、位置による分布

の偏りが少なくなり、それだけ重要部と不要部の区別が付きにくくなる。よって、「切捨て」効果もそれだけ薄れるのではないかと推察される。

また、PLMの問題点はテキストのサイズが小さく、またブロックサイズ(j 値)が小さいと、再現率(recall)が十分に伸びず、break-even point が得られないことである。理由は、ブロックに含まれる単語が少なすぎるという点にあると考えられる。表3で $j = 20\%$ のところをみると、精度が全文方式に比べて17%も落ちているが、 j を大きくしていくにつれパフォーマンスは向上する。興味深いのは、 l が大きく j が小さい場合である。例えば、 $j = 20$ のとき、対フルテキストで $500 < l \leq 600$ で20%の精度アップ、反対に $100 < l \leq 200$ のとき17%のダウンが見られる。

最後にFPM(first-paragraph model)について結果を述べる。(表1参照。) FPMは記事を第一段落のみを用いてその記事の見出し語を推定するアプローチである。結果的には、パフォーマンスにむらがあり、他方式に比べよい場合(300-400, 500-600)もあるが極端に落ちることもある(900-1000)。特に、900文字以上のテキストについては、全文方式よりも悪い。図6に3つの方式をグラフにしたものを載せた。結果的には図をみると、FLM方式で $i = 20$ のとき、精度が最もよかつ全体的に安定していると思われる。

7 まとめ

主題推定にテキスト構造を利用する最大のメリットは、精度の向上と推定に必要なデータの圧縮である。後者については、FLM方式で20ブロックの場合ではフルテキストの場合に比べて必要なデータは10分の1に圧縮されている。本稿では、主題推定の基本的な方式を述べ、この方式が長文テキストに対して弱いという問題点を提示した。問題の解決法として、本稿では文章の重要部のみを用い、不要部を主題推定の対象から除外するという方法を採用した。結果的には、全体的に精度の向上が確認された。特に、テキスト記事の長さに関係

表 1: Results for using paragraph

Size	100-200	200-300	300-400	400-500	500-600	600-700	700-800	800-900	900-1000
Break Even	.396	.358	0.389	.371	.381	.338	.290	.283	.250

表 2: Fixed-Length Model (FLM). i はブロックの長さ、 l はテストセット。 $i = 10$ のカラムは、記事の先頭 10 単語をから記事の主題 (見出し語) を推定した場合の結果である。長さ 400 文字以上の記事では break-even point が存在しない。

$l \setminus i$	10	20	30	40	50	60	70	80	90	100
100-200	.38(-.10)	.42(-.00)	.41	.41(-.02)	.41	.41(-.02)	.41	.41(-.02)	.41	.41(-.02)
200-300	.35(+.03)	.39(+.15)	.38	.36(+.06)	.35	.34(+.00)	.34	.34(+.00)	.34	.34(+.00)
300-400	.34(+.06)	.38(+.19)	.37	.36(+.12)	.35	.34(+.06)	.33	.33(+.03)	.33	.33(+.03)
400-500	*.34(.30)	.37(+.19)	.37	.36(+.16)	.35	.34(+.10)	.33	.33(+.06)	.32	.32(+.03)
500-600	*.34(.28)	.37(+.23)	.36	.36(+.20)	.36	.35(+.17)	.33	.34(+.13)	.33	.33(+.10)
600-700	*.36(.29)	.37(+.32)	.37	.36(+.29)	.35	.33(+.18)	.33	.32(+.14)	.31	.31(+.11)
700-800	*.32(.25)	.34(+.31)	.34	.34(+.31)	.33	.32(+.23)	.31	.30(+.15)	.30	.29(+.12)
800-900	*.32(.25)	.34(+.31)	.34	.33(+.27)	.33	.32(+.23)	.32	.31(+.19)	.30	.30(+.15)
900-1000	*.29(.23)	.32(+.23)	.31	.30(+.15)	.30	.29(+.12)	.29	.29(+.12)	.28	.28(+.08)

なくテキストの始まりから特定の長さのブロックを切り出し、推定に用いた場合に、最もよい結果が得られた。また、テキストの第一パラグラフを使って同様の実験を行なったが、われわれの予想に反し、あまり効果が確認できなかった。今後の予定としては、本研究の知見をもとにスキミング (飛ばし読み)、スキニング (探し読み) などの人間の読書能力を効果的に支援する手法について考えていきたい。

参考文献

- (Callan, 1994) James P. Callan. Passage-Level Evidence in Document Retrieval. In Croft and van Rijsbergen (1994), pages 302–310.
- (Croft and van Rijsbergen, 1994) W. Bruce Croft and C. J. van Rijsbergen, editors. *Proceedings of the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. Dublin City University, Springer-Verlag, 1994.
- (Fox, 1987) Barbara A. Fox. *Discourse structure and anaphora*. Cambridge Studies in Linguistics 48. Cambridge University Press, Cambridge, UK, 1987.
- (Hearst, 1994) Marti A. Hearst. Multi-Paragraph Segmentation of Expository Text. In *Proceedings of the 32nd Annual Meeting of the Association for Computational Linguistics*, pages 9–16, New Mexico, USA, 1994.
- (Lewis, 1992) David D. Lewis. An Evaluation of Phrasal and Clustered Representations on a Text Categorization Task. In *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 37–50, 1992.
- (Matsumoto *et al.*, 1994) Yuji Matsumoto, Sadao Kurohashi, Takehito Utsuro, Yutaka Myoki, and Makoto Nagao. Japanese Morphological Analysis System JUMAN Manual. NAIST-IS-TR 94025, Nara Institute of Technology and Science, Nara, Japan, 1994.
- (Salton *et al.*, 1993) Gerald Salton, J. Allan, and C. Buckley. Approaches to Passage Retrieval in Full Text Information Systems. pages 49–58. ACM Press, 1993.
- (Wilkinson, 1994) Ross Wilkinson. Effective Retrieval of Structured Documents. In Croft and van Rijs-

表 3: Proportional-Length Model (PLM). j はブロックの記事に対する割合、 l はテストセット。

$l \setminus j$	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
100-200	*.40(.22)	.35(-.17)	.39	.41(-.02)	.42	.42(-.00)	.41	.41(-.02)	.41	.42
200-300	*.38(.24)	.37(+.09)	.39	.38(+.12)	.36	.36(+.06)	.36	.35(+.03)	.35	.34
300-400	*.36(.30)	.38(+.19)	.38	.36(+.12)	.36	.35(+.09)	.34	.34(+.06)	.33	.32
400-500	*.34(.33)	.38(+.23)	.37	.35(+.13)	.34	.34(+.10)	.33	.33(+.06)	.32	.31
500-600	.34(+.13)	.36(+.20)	.36	.35(+.17)	.34	.33(+.10)	.33	.32(+.07)	.32	.30
600-700	.35(+.25)	.36(+.29)	.35	.33(+.18)	.32	.31(+.11)	.30	.29(+.04)	.29	.28
700-800	.36(+.38)	.34(+.31)	.32	.32(+.23)	.29	.28(+.08)	.28	.27(+.04)	.26	.26
800-900	.35(+.35)	.33(+.27)	.32	.31(+.19)	.30	.29(+.12)	.28	.27(+.04)	.27	.26
900-1000	.32(+.23)	.30(+.15)	.29	.28(+.08)	.27	.27(+.04)	.26	.26(+.00)	.26	.26

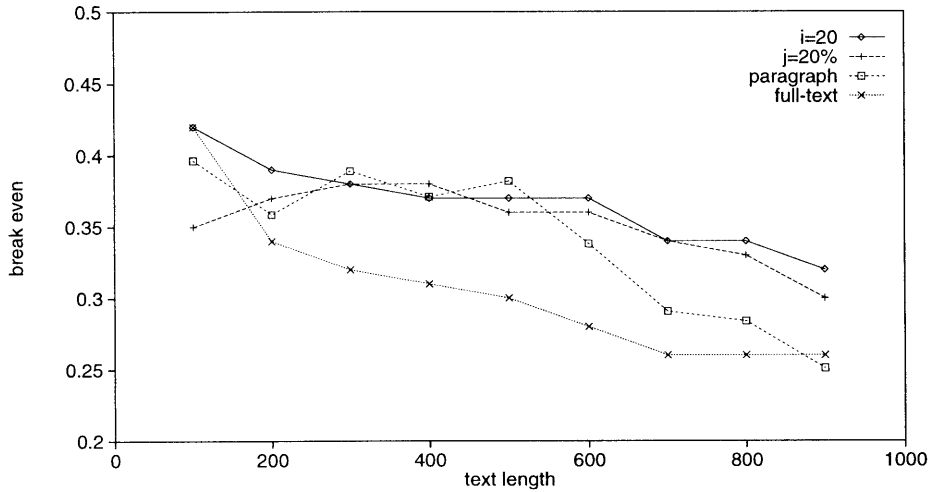


図 6: FLM($i = 20$), PLM($j = 20\%$), 「第一段落」, 全文, 各方式の比較

bergen (1994), pages 311-317.

(日本経済新聞社, 1992) 日本経済新聞社. 日本経済新聞
92年 CD-ROM 版, 1992. 東京.