

ニュース映像のシーン順序に基づく Blog 検索

北山 大輔[†] 角谷 和俊^{††}

[†] 姫路工業大学環境人間学部 〒 670-0092 兵庫県姫路市新在家本町 1 丁目 1-12

^{††} 兵庫県立大学環境人間学部 〒 670-0092 兵庫県姫路市新在家本町 1 丁目 1-12

E-mail: [†]na01a066@stshse.u-hyogo.ac.jp, ^{††}sumiya@shse.u-hyogo.ac.jp

あらまし ニュース映像には、いくつかのシーンが含まれ、ニュースの構成要素となっている。ユーザはこれらのシーンに基づき、そのニュースの概要を理解することが可能である。しかし、シーンの出現の順序によって、ユーザのニュースに対する評価が変化する場合がある。すなわち、一つのニュースに対して、複数の評価が存在することが考えられる。本研究では、ニュース映像サイトにおいて、シーン順序を自動構成する環境を想定し、シーンの出現パターンによる評価の差異を利用し、その評価と同じ Blog 記事を抽出する方式を提案する。本稿では、シーン順序に基づくニュースの評価手法について述べ、ニュース映像と Blog 記事を用いた情報融合環境について議論する。

キーワード ニュース映像, シーン順序, Blog 検索, 評価抽出

Blog Search using News Video Scene Order

Daisuke KITAYAMA[†] and Kazutoshi SUMIYA^{††}

[†] School of Humanities for Environmental Policy and Tecnology, Himeji Institute of Technology

1-1-12 Shinzaike-honcho, Himeji, Hyogo 670-0092, Japan

^{††} School of Human Science and Environment, University of Hyogo

1-1-12 Shinzaike-honcho, Himeji, Hyogo 670-0092, Japan

E-mail: [†]na01a066@stshse.u-hyogo.ac.jp, ^{††}sumiya@shse.u-hyogo.ac.jp

Abstract News video stream consists of several scenes. Users can realize news outline by viewing some scenes. However it sometimes happens that they change the news evaluation by scene order. In other words, we consider that one news video has many evaluations. We propose a method to extract Blog entries that have same evaluations using the difference of evaluations by news sequence. In this paper, we describe the method of news evaluation based news scene order. We also discuss information integrated environments using Blog and news video stream.

Key words news video, scene order, Blog search, evaluate extraction

1. はじめに

映像メディアによるニュースは重要な情報源となっている。近年、ニュース映像は TV のみならず FNN-NEWS.COM [1] などのインターネット上のニュースサイト [2] ~ [4] でも配信されるようになってきた。しかしながらその利用は、従来の映像メディアと同じ方法で展開されている。従来の映像メディアでは時系列的に作成者の意図したとおりの見方しかすることができないが、意見の対立や立場の違いなど複数評価を持つニュース映像では、シーンの構成順序によって、映像から受ける評価の印象が変化すると考えられる。例えば、「鯨肉の給食が再開される」という捕鯨問題に関するニュースがあるとすると、前半シーンで、「おいしい」のように歓迎的な意見が述べられていても、後半シーンに「今の日本に鯨を食べる習慣は無い」といった否

定的な意見の内容があれば、あまり歓迎されていない印象を受ける。反対に、前半で否定されていても後半で歓迎の意見があれば歓迎されている印象を受ける。このようにシーン順序が異なると、ユーザが受ける評価の印象に違いが生じると考えられる。

一方、一人の書き手の物でしかなかった Web ページだが、Blog, BBS といった複数の書き手が意見を交換できる環境が整ってきた。その中の一つで書き手が積極的に意見を書き、その意見に対して批評を行なうことができる Blog が広まってきている。Blog の中にはニュースの事柄に対して評価を行なっているものも少なくない。複数評価に分かれるようなものは、Blog で議論が交わされることがあり、同じトピックであっても議論の肯定・否定の割合が異なる Blog がある。

ニュース映像では、時間的な制約から各々の評価に対して十

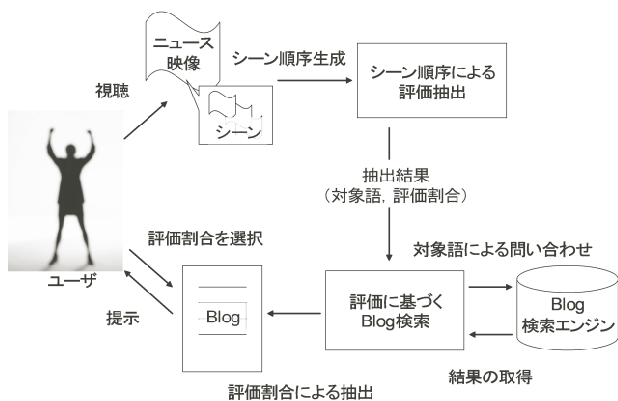


図1 評価抽出に基づく Blog 検索の概念図

分な解説がなされないことが多い。そこでニュース映像のシーン順序によりニュースの評価の割合を抽出し、その評価割合に合う Blog を検索する手法を提案する。ニュース映像から得られた評価割合と合う Blog 記事を得ることで、映像ニュースでは不足している様々な角度からの評価を得ることができる。

以下、本稿の構成を示す。まず2節では本研究の概要と関連研究について説明する。3節ではシーン順序による評価の抽出方法について説明する。4節では評価を用いた Blog 検索について説明する。5節では作成したプロトタイプシステムについて説明し、評価実験とその考察を行なう。最後に6節でまとめと今後の課題について述べる。

2. 本研究の概要と関連研究

2.1 本研究の概要

本研究では、ニュース映像のシーン順序によりニュースの評価の割合を抽出し、その評価割合に合う Blog を検索する手法を提案する。図1は本手法の概念図を示す。本手法の特徴は以下のとおりである。それぞれ図1の上部と下部に該当する。

- ニュースの取りうる評価を自動で抽出する。

ニュースのシーン順序により、ニュースの取りうる評価が変化すると考え、シーン順序の生成とそのシーン順序に基づく評価の抽出を行なう。シーン順序生成は順序変化により影響を与えやすいシーンを抽出し、自動で順序を生成する。生成された順序毎に対象語・評価表現の重要度計算を行ない、ニュースの評価として対象語と評価割合を抽出する。

- 評価に基づく Blog の提示を行なう。

Blog に評価が含まれる場合、同一トピックに対する Blog であってもさまざまな評価を行なっていると考え、キーワードによる検索の後、評価表現の重要度算出を行ない、評価割合によってクラスタリングし、評価割合ごとに提示する。

上記の2つにより、ニュース中のトピックについて評価を行なっている Blog をニュース中の評価割合に基づき検索を行なうことが出来る。本手法によりユーザはニュースを見ることで、そのトピックについての様々な評価を得ることが出来る。

2.2 関連研究

ma ら [5] や Henzinger ら [6] は映像の音声テキストを用いて

自動で検索質問を作成し、Web より情報を取得する手法を提案している。音声テキストから検索質問のキーワードを抽出し検索を行なうという点で本研究と類似しているが、ニュース映像中の評価に基づき、Blog の検索を行なうという点で目的が異なる。

若林ら [7] は個人の嗜好と評判情報を用いて商品購入予定者に最適な評判情報を提示する手法を提案している。購入経験者の嗜好の偏りと購入予定者の嗜好の偏りのマッチングを行ない、該当する評判情報を提示することで、ユーザにとって有益な情報を提示することができる。評価をマッチングし情報を提示するという点が同じであるが、対象が商品情報ではなくニュース映像中のトピックである点と、ユーザに対象に対する多面的な見方を提示することを目的とする点で異なる。

奥村ら [8] は blogWatcher という Blog 検索に特化したシステムを構築している。blogWatcher では、Blog 検索以外に、キーワードの burst 度、ホットキーワードの表示、評判情報検索の機能を有している。このうち、評判情報検索は入力キーワードに関する評判情報を Blog から検索して提示するという部分が本研究の Blog 検索と似ているが、Blog 中における議論の肯定・否定の割合に着目し、評価割合の一致により検索を行なっている点で異なる。

Kumuar ら [10] はハイパーリンクによる Blog 群のつながりを blogspace と定義し、Blog コミュニティの抽出を行ない、コミュニティの進化に関する研究を行なっている。Gruhl ら [11] はマクロな視点、ミクロな視点からのトピック伝達の特徴づけを行ない、議論が行なわれている Blog と外部的要因により発生する Blog により Blog における情報伝播のモデル化を行なっている。ただしこれらの研究は、情報の伝播に着目したものであり Blog 中のトピックに対する評価を扱うものではない。

Blog 検索に関する先行研究として、竹原ら [12] や中島ら [13], [14] の研究が挙げられる。いずれも Web ページの信頼性を算出する手法として Blog を用いており、評価に基づく Blog 検索を行なう本研究とは、手法も目的も異なる。また、是津ら [15] の研究が挙げられるが、コンテンツそのものを対象とした閲覧のための手法であり、本研究とは対象が異なる。

評価表現の抽出および評価判別に関する研究として、藤村ら [16]~[18] や鈴木ら [19]、立石ら [20] の研究があげられる。本研究では、評価表現辞書の作成に関して、キーワードベースで評価の判別を行なう立石らのアプローチを参考にした。

2.3 予備実験

2.3.1 シーン順序による評価の生成

報道各局により同一トピックのニュースであっても、その評価が異なる場合が存在する。例えば、郵政民営化のニュースについて、TBS では小泉総理が苦言を呈するシーンを最後に使い、自民党からは不満の声が出ているというように否定寄りにニュースを構成している。一方、FNN では小泉総理の苦言シーンを前半部に使い、自民党に関しては内容を協議するといった表現にとどめている。

TBS のニュースのシーン順序を入れ替えて、FNN のニュースと類似したシーン順序を作ることで、FNN に近い評価を生成

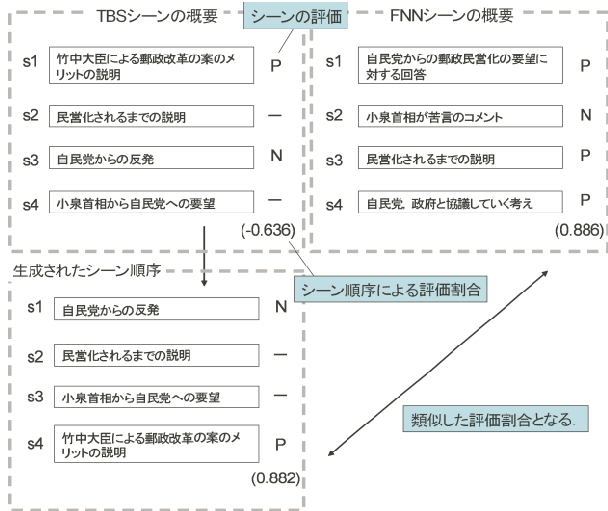


図2 予備実験

することが可能であるかどうかの予備実験を行なった。手順は以下の通りである。

1. FNN のニュース映像におけるシーンの評価付けを行なう。
2. TBS のニュース映像におけるシーンの評価付けを行なう。
3. TBS のニュース映像に含まれるシーンの全組み合わせを自動生成する。
4. 自動生成されたシーン順序から手作業で FNN と近いシーン順序のものを抽出する。

2.3.2 結果と考察

上記の実験に挙げた例では、TBS のシーンを入れ替えることで FNN に近い評価を生成することができることを確認した。図2に実験の様子を示した。図中の文章は各シーンの概要、文書の右にある記号は P が肯定シーン、N が否定シーン、- が評価のないシーンを示している。FNN では結論となる最後の部分に肯定シーンが連続し、肯定的な印象を受ける順序となっている。そのため、TBS では肯定シーンを結論にあたる最後にし、肯定の印象を強めるために否定シーンを最初とする順序がもっとも FNN に近い印象と考えられる^(注1)。以上より、シーン順序を入れ替えることで他の局のニュースの評価の印象に近い順序を作り出せることを確認した。

3. シーン順序による評価の抽出

3.1 シーンの評価と操作

ニュース映像のシーン分割には、時間による分割、話題の区切れによる分割、画面切り替えによる分割など、さまざまな方法があるが、本稿では画像処理や音声認識などを用いて分割する。音声の無音区間を軸とし、その周辺に映像の切り替えが存在するならば、その無音区間をシーンの切れ目として扱う。なお、シーンの操作に関してはユーザが自由に行なうことも考えられるが、本稿ではシステムが自動で順序生成を行なう。

(注1): 3節で提案する評価割合の算出方法で計算を行なったところ、ほぼ同様の値が得られた。

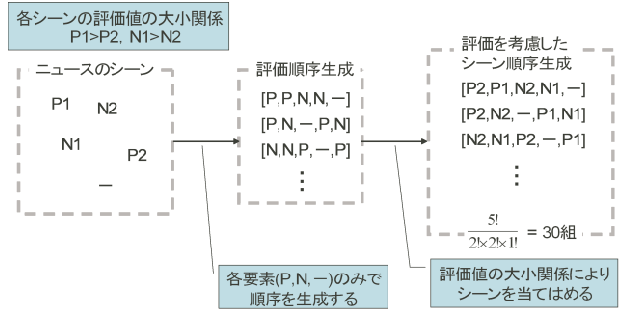


図3 評価を考慮した組み合わせの削減

3.2 評価表現

評価表現を判別するために評価表現辞書を作成した。辞書には“歓迎”、“批判”のような肯定もしくは否定の意味を含む単語を肯定・否定という情報とともに登録している。また、株価の“上昇”、“下落”など、一見事実であるように見えるが、評価とも取れるものは評価表現として辞書に登録している。この辞書に基づいて評価の肯定・否定の判別を行なった。また、評価の補助表現として“ない”のように評価表現と組み合わせることで評価が反転するような単語も登録している。評価表現の近傍に補助表現が存在する場合は簡易な言語処理により評価の肯定・否定を反転させている。なお、疑問文中に出現する評価表現には評価の意図はないと考えられるが、今回は考慮していない。

3.3 組み合わせ数の削減

シーン順序生成を自動で行なう場合、総当りで行なうと組み合わせ数が多くなるという問題がある。シーン数の多い映像では、膨大な計算を必要とするため、シーン操作の組み合わせ数を適切に削減する方法が必要となる。

シーンの肯定、否定という評価に基づいて組み合わせの削減を行なう^(注2)。シーン s_i について、シーン中に含まれる評価表現 e_{ij} をシーン内の出現順により重要度算出を行なう。算出された肯定重要度と否定重要度を減算することによってシーン評価値 $S_{eval}(s_i)$ を決定する。式中の E_i はシーン内の評価表現出現数である。

$$S_{eval}(s_i) = pos(e_{ij}) - neg(e_{ij}) \quad (1)$$

$$pos(e_{ij}) = \sum_{j=1}^{E_i} (e^{-(E_i-j+1)}) \quad (2)$$

$$neg(e_{ij}) = \sum_{j=1}^{E_i} (e^{-(E_i-j+1)}) \quad (3)$$

なお、 $pos(e_{ij})$ では e_{ij} が肯定である場合のみ重要度の算出を行ない、 $neg(e_{ij})$ では否定の場合のみ重要度の算出を行なう。 $S_{eval}(s_i)$ が正の値をとればシーン評価を肯定、負の値をとれば否定とする。評価表現の含まれないシーンの場合 $S_{eval}(s_i)$ がゼロとなる。そのような場合にはシーン評価を評価無しとした。

まず、シーン評価の否定、肯定、評価無しの3要素から作ら

(注2): 評価表現のカテゴリ、評価表現の指し示す対象に関しては本稿では考慮しない。

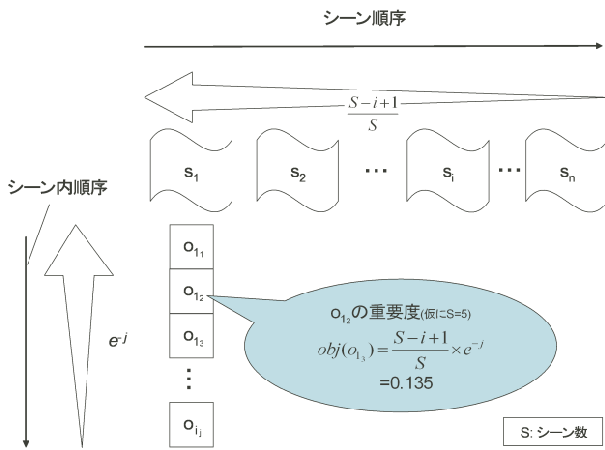


図4 対象語の重要度の算出

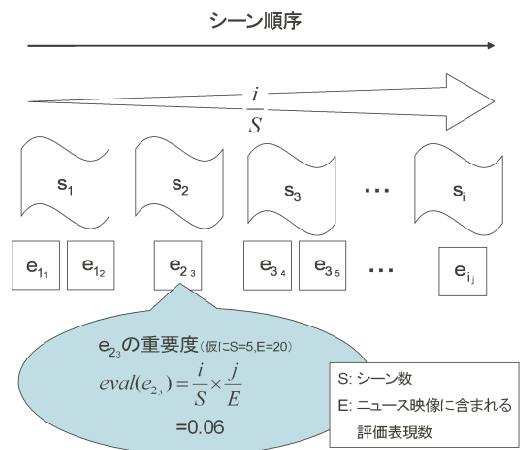


図5 評価表現の重要度の算出

れうるすべての組み合わせを作成する．その後シーン評価が同じものに関して、評価値の絶対値が大きいものほど後のシーンとすることで組み合わせの削減を行なう．評価値が大きいシーンを後とするのは、評価に影響を与えやすくするためである．逆に言うと、評価に影響を与えないような組み合わせの削減を行なっている．評価無しシーンに関しては順序は基となるニュース映像のシーン順序とした．総当りであればシーン数の階乗個生成される組み合わせを、この手法によりシーン数の階乗を肯定シーンの階乗、否定シーンの階乗、評価無しシーンの階乗の3つを乗算した値で除算した数の組み合わせに削減することができる．

肯定シーン P が2つ、否定シーン N が2つ、評価無しのシーン $-$ が1つの時の例を図3に示した．各シーンの評価値の大小関係は $P_1 > P_2$, $N_1 > N_2$ である．この例では、総当りの組み合わせが120組生成されるのに対し、30組に削減できる．

3.4 ニュースのトピックに対する評価

3.4.1 対象語の抽出

対象語とは、評価の対象となるトピックを表す名詞である．ただし、名詞のうち対象となりえない品詞、例えば代名詞などは除く．ニュース映像においては、シーン順序が先の方ほど重要であり、シーン内でも先に出現するキーワードが対象語として重要であるといえる．また、ニュース内で複数回出現するようなキーワードも対象語として重要といえる．そのため、シーン順序が先であるほど重要さが上がり、シーン内での出現順が先であるほど重要さが上がるような関数を用いて抽出を行なう．図4に算出の手順を示す．

シーン中の名詞 o_{ij} について、シーン順重要度とキーワード出現順重要度を乗算し、重要度 $obj(o_{ij})$ を算出する．同一キーワードに関しては、同一キーワード同士の値を足し合わせたものを対象語の重要度とする．式中の i はシーン順序、 j はシーン内の出現順であり S はシーン数である．

$$obj(o_{ij}) = \frac{S-i+1}{S} \times e^{-j} \quad (4)$$

名詞に関しては茶筌[21]による形態素解析を用いて抽出を行なう．対象語の重要度とは、そのニュース映像内で評価対象と

して扱われうる候補となるかどうかを表す値で、大きければニュース映像内で評価の対象となっている可能性が高い．この対象語の重要度計算を各シーン順序毎に行なうことで各シーン順序の対象語 O を抽出する．対象語 O は、各シーン順序の重要度上位の名詞の集合である．

3.4.2 評価割合の算出

本節では、評価表現の抽出方法について述べる．評価とは、評価の対象を表す対象語とその評価内容を表す評価表現の組み合わせである．

評価表現は対象語とは異なり、ニュース映像内での出現順序が後であるような評価表現ほどユーザに対して強い印象を与えるという特徴がある．例えば「環境保護団体は鯨保護のために捕鯨に反対する」と「環境保護団体が捕鯨に反対なのは鯨保護の精神からである」であれば、前者のほうが反対という印象が強いといえる．そのため、シーン順序が後であるほど重要さが上がり、ニュース映像内のすべてのキーワードの中での評価表現の出現位置が後であるほど重要さが上がるような関数を用いて重要度の算出を行なう．

ニュース映像中の評価表現 e_{ij} について、出現シーン順とニュース映像内出現順を用いて重要度 $eval(e_{ij})$ を算出する．評価表現は対象語とは違い、同一キーワードでも評価が反転する場合が考えられるため、同一キーワードの足し合わせを行なわない．図5に算出の手順を示す．式中の i はシーン順、 j はニュース内の評価表現の出現順、 S はシーン数、 E はニュース映像内の評価表現の総数である．

$$eval(e_{ij}) = \frac{i}{S} \times \frac{j}{E} \quad (5)$$

評価表現に関しては評価表現辞書を用いて抽出を行なう．評価表現の重要度とは、そのニュース映像内でいかにユーザに印象を与えるかどうかを表す値で、大きければニュース映像内でユーザにその評価の印象を与えている可能性が高い．出現するすべての評価表現中、肯定表現の重要度を足し合わせた値を P_{video} 、同じく否定表現の重要度を足し合わせた値を N_{video} とし、以下の式でシーン順序の評価割合 V_ratio を算出する．

$$V_ratio = \frac{P_{video} - N_{video}}{P_{video} + N_{video}} \quad (6)$$

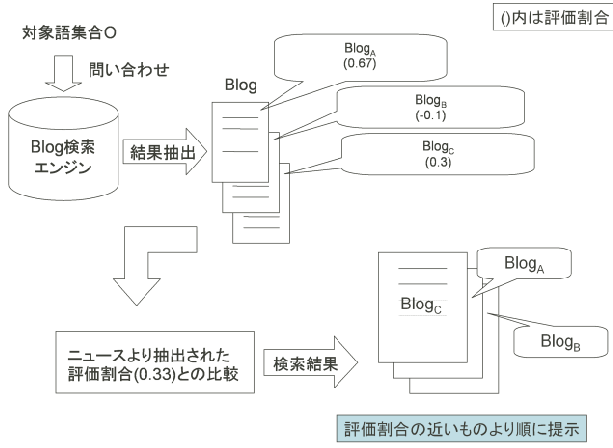


図6 評価に基づく Blog 検索

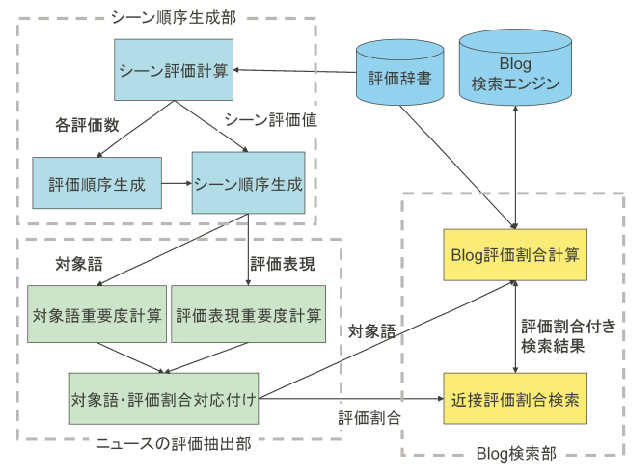


図7 システム構成図

$$P_{video} = \sum eval(e_{i_j}) \quad (7)$$

$$N_{video} = \sum eval(e_{i_j}) \quad (8)$$

なお、 P_{video} では e_{i_j} が肯定的場合に算出を行ない、 N_{video} では否定の場合に算出を行なう。評価割合の取りうる値は -1 から 1 の間であり、 -1 に近いほど否定寄り、 1 に近いほど肯定寄りのシーン順序の組み合わせとする。

4. 評価を用いた Blog 検索

4.1 Blog 検索の目的

通常の Blog 検索では、検索キーワードが Blog のエントリ中に含まれているような Blog を探し出して提示する。ニュース映像の対象とその評価を議論する Blog を検索するには従来のキーワードベースの手法では限界がある。

評価を用いた Blog 検索の目的とするところは、ニュース映像から抽出された評価対象について、抽出される評価の偏り具合を表しているような Blog を結果として提示することにある。例えば、捕鯨のニュースに関して、捕鯨の給食について否定寄りの印象を受けるシーン順であれば、取得される Blog は、エントリで捕鯨給食について扱っており、かつコメント・トラックバックを含めた議論が否定寄りで展開されていることが望ましい。

Blog はエントリタイトル、エントリ、コメント、トラックバックから構成されることとし、それぞれを区別して扱う。検索は、対象語を用いてエントリタイトル、エントリに対して行なう。これは、コメント中のみ対象語が現れる Blog はそのトピックについて扱っていないと考えられるからである。Blog の評価割合を算出する際には、エントリ、コメント、トラックバックの評価表現の重要度より算出を行なう。別の Blog であるトラックバック先からも評価を抽出しているのは、コメントと同様にエントリから派生した評価と考えられるからである。

4.2 Blog に含まれる評価表現の重要度

本節では、Blog に含まれる評価表現の重要度について述べる。Blog では文章を構成するときに文章の後半に自分の意見を述べることが多いため、書き手の意見として文書の後半の表現ほど重要であるといえる。そのため、エントリ、個々のコメン

ト、トラックバックに対して評価表現辞書を用いて評価表現の抽出を行ない、文章中の出現順序が後の方であるほど重要さが上がるような関数を用いてエントリ、個々のコメント、トラックバックに含まれる評価表現の重要度の算出を行なう。式中の B は算出対象となるエントリ、個々のコメント、トラックバックに含まれる評価表現の個数である。

$$blog(e_i) = e^{-(B-i+1)} \quad (9)$$

Blog 中の評価表現の重要度は書き手の意見らしさを表しており、重要度の高い評価表現ほど書き手の意見といえる。

4.3 Blog の検索

映像より抽出された対象語 O と評価割合 V_{ratio} を用いて検索を行なう。Blog 検索エンジンに対し、対象語を検索キーワードとして検索を行なう。得られた Blog 集合中、エントリタイトルおよびエントリ中に対象語が現れているものを対象 Blog 集合とする。個別の対象 Blog 中のエントリ、コメント、トラックバックより評価表現を抽出し、重要度により評価表現の割合を算出する。Blog 中の肯定表現の重要度を足し合わせた値を P_{blog} 、否定表現の重要度を足し合わせた値を N_{blog} とし以下の式で Blog の評価割合 B_{ratio} を算出する。

$$B_{ratio} = \frac{P_{blog} - N_{blog}}{P_{blog} + N_{blog}} \quad (10)$$

$$P_{blog} = \sum blog(e_i) \quad (11)$$

$$N_{blog} = \sum blog(e_i) \quad (12)$$

シーン順序の評価割合と対象 Blog の評価割合を比較し、割合の近いものを抽出し提示する。これをシーン順序により抽出される評価割合すべてに対して行なう。この結果、あるニュースの映像のシーンの組み合わせから自動的に生成されたいくつかの評価割合に基づき、対応する Blog が抽出される。

5. プロトタイプシステム

5.1 システム構成

プロトタイプシステムは大きく分けてシーン順序生成部、評価抽出部、Blog 検索部の3つからなる。図7にシステム構成を

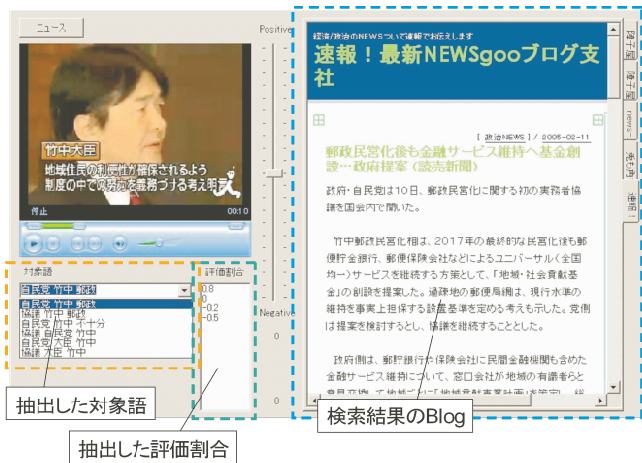


図8 プロトタイプシステム

示した。

シーン順序生成部、評価抽出部によりニュース中の対象語と取りうる評価割合の抽出を行なう^(注3)。Blog 検索部において、取得した Blog 中の評価割合をクラスタリングし、ユーザへの提示を行なっている^(注4)。なお、シーン分割部分を実装していないため、シーンの分割結果は映像分割、音声分割ともに他のシステムにより分割し、その結果を用いて人手で生成したものを使用している。分割の方法は以下のとおりである。

- 無音区間の前後 1 秒に映像による区切れがあればシーンの区切れとする^(注5)。
- 一つの映像による区切れが複数の無音区間の前後 1 秒に含まれる場合は、もっとも近い無音区間をシーンの区切れとする。

5.2 実装環境

プロトタイプで用いたニュース映像は FNN-NEWS.COM、TBS NEWS i などの映像ニュースサイトより取得した実際のニュース映像であり、Blog の情報は goo ブログ [22] の検索を用いて取得した実際の Blog 記事である。ニュース映像のシーン組み合わせの生成および対象語、評価表現の重要度計算、Blog 記事の評価割合の計算部分には Visual Studio.NET の C# を用いた。ニュース映像の音声テキストおよび Blog 記事中の単語抽出には茶釜による形態素解析を用いて抽出した。評価表現の抽出には、人手により作成した評価表現辞書を用いている。実験に用いたマシンは OS:WindowsXP Home/ CPU:Celeron 2.53GHz/ メモリ:512MB である。

5.3 ユーザインタフェース

図8は実際のプロトタイプの画面である。画面左上の部分でニュース映像を視聴することが出来る。システムはニュースのシーンの音声テキストを用いて、対象語と評価割合を抽出、対

(注3): 「ない」などの評価補助語による評価の反転は最大で一回のみとし、肯定+二重否定という場合には対応していない。

(注4): Blog 中のトラックバックは利用していない。

(注5): 自動検出の無音区間の前後 1 秒というのは音声と映像のずれの誤差を考慮したためである。

象語により Blog 検索エンジンに問い合わせ、検索結果の Blog の評価割合を計算する。以降ユーザが評価の知りたい対象語を選択すれば、その対象に関する Blog を検索結果として提示する。また、ニュース映像から抽出された評価割合を選択することで、ニュースの評価割合と Blog の評価割合を比較し、割合の近いものを検索結果としてユーザに提示する。

5.4 システム処理の流れ

シーン順序生成部の処理の流れは以下の通りである。

1. ユーザにより映像ファイルが選択される。
2. 映像ファイルを基にあらかじめ作成しておいた各シーンの音声テキストをシステムが読み込む。
3. 音声テキストを形態素解析により単語に分割する。
4. 形態素解析結果より名詞と出現順を抽出し対象語を得る。同じく評価表現も出現順と共に抽出する。
5. 評価表現に関し出現順に基づく重要度の算出を行なう。
6. 算出結果を用いシーンの否定肯定評価を決定する。
7. 肯定、否定、評価無しそれぞれのシーン数を元にシーンの組み合わせを生成する。
8. 否定・肯定より生成された組み合わせにシーン評価値を基にシーンの順序を生成する。

評価抽出部の処理の流れは以下の通りである。

1. 対象語にシーン順序を付与し、重要度の計算を行なう。
2. 対象語の重要度上位 3 個をそのシーン順序の対象語とする。
3. 評価表現にニュース内出現順を付与し、重要度の算出を行なう。
4. 評価表現の否定・肯定ごとに重要度を足し合わせ評価割合を算出する。

Blog 検索部の処理の流れは以下の通りである。

1. 抽出された対象語を検索キーワードとして、goo ブログにより AND 条件で検索を行ない結果上位 100 件までを取得する。
2. 取得した Blog 中より HTML 解析により、エントリとコメントのみを抽出する。
3. エントリ、コメント中より評価表現を抽出し、重要度の算出を行なう。
4. 重要度の算出結果を用いて、否定・肯定の重要度を足し合わせ、評価割合の算出を行なう。
5. シーン順序の評価割合と Blog の評価割合を比較し、割合の近いものを 5 件程度ユーザに提示する。

シーン順序生成部により生成されたシーン順すべてに対して評価抽出部を実行し、抽出された対象語の組み合わせすべてに対し検索を行なう。

5.5 評価実験

5.5.1 実験方法

対象語抽出に関する評価をするため、出現頻度により抽出されたキーワードとの検索結果の比較実験を行なった。評価実

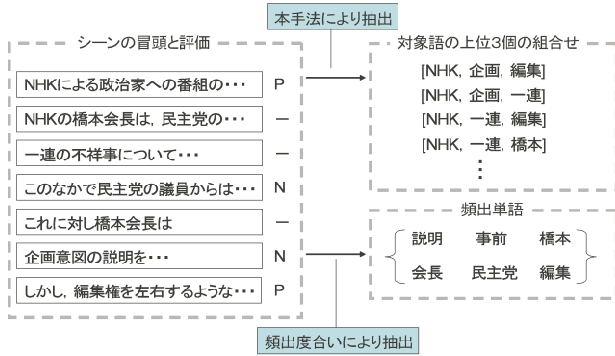


図9 評価実験のデータ

験に用いたデータを図9に示す。右側は各シーンの文頭とそのシーンの評価であり、左下部は対象語により生成されたキーワードの組み合わせ、左上部は頻出度合いの上位となったキーワードである。実験の際に用いた Blog データは、ニュースの前後1週間の goo ブログの検索結果を取得したものである。

対象語による質問は、シーン順により出現する対象語上位3位を用いて生成し、AND 条件で検索する。シーン順序は複数生成されるので、そのすべてに関して質問を生成し、適合率を計算する。頻出単語による質問は、頻出度順上位3位までのキーワードから作られる組み合わせすべてにより生成し、その組み合わせごとに AND 条件で検索を行なう。また、それぞれの組み合わせに対して適合率を計算する。なお、頻出度順第3位に複数のキーワードが存在する場合、そのすべてを用いて組み合わせを生成する。

全体集合、正解集合を共通とし、全体集合は対象語から生成される質問に含まれるキーワード、頻出単語による抽出で出現するキーワードのすべてを OR 条件で検索した結果とする。正解集合の異なる2つの実験を行なった。

(実験1) 正解集合を全体集合より手作業で抽出した「対象としたニュースのトピックについて触れている」ような Blog の集合とした実験。

(実験2) 正解集合を全体集合より手作業で抽出した「対象としたニュースのトピックについて何らかの評価を行っている」ような Blog の集合とした実験。実験2の集合は、実験1の集合に含まれる。

この2つの実験の適合率を比較して評価を行なう。実験の手順を以下に記す。

1. 頻出単語および対象語の抽出を行なう。
2. goo ブログより頻出単語、対象語一つ一つに対して検索を行ない一週間以内のものを取得し全体集合とする。
3. 手作業により実験1, 2の正解集合を作成する。
4. それぞれのキーワード組を AND 条件で検索をかけ、結果を取得する。
5. 検索結果と正解集合より、適合率の計算を行なう。

5.5.2 結果と考察

実験データとして、NHK 会長が番組改変問題の不祥事について謝罪と説明を行なうニュースを用いた。ニュースの内容より、

表1 評価実験の適合率

		実験1	実験2
対象語	適合率	0.875	0.542
頻出単語	適合率	0.221	0.097

正解集合を番組改変問題に触れている Blog として実験を行なった。対象語にのみ現れたキーワードは“NHK”、“企画”、“一連”、頻出単語として抽出されたものは“説明”、“事前”、“会長”、双方に共通して現れたキーワードは“橋本”、“民主党”、“編集”というキーワードであった。対象語として7組の質問が生成され、頻出単語からは20組を生成した。そのうち検索結果が得られたものは、対象語では4組、頻出単語では11組であった。実験として、検索結果の返ってきた組み合わせを比較した。表1は実験1, 2における対象語、頻出単語それぞれの適合率の平均値である。以下に考察を述べる。

- 頻出単語との比較を行なうと、全体的に対象語の方が高い値をとっていることがわかる。頻出単語では、検索の解に正解を含まない組み合わせが多数存在するのに対し、対象語では検索の解に正解を含まない解がなかった。このことより対象語は正解に対し、的確に質問を生成していると考えられる。

- 実験1の結果と実験2の結果の比較を行なうと、頻出単語では値が大幅に減少するのに対し、対象語では頻出単語ほどの減少は見られない。実験1と実験2の違いは、正解中に評価を含むかどうかである。このことより、頻出単語により抽出された解は評価をあまり含まないのに対し、対象語により抽出された解には評価を含む解が多いことが考えられる。

- 検索結果の Blog の内容についての比較を行なうと、頻出単語ではニュースそのものを扱った Blog が多く見られた。今回の実験では NHK 会長が番組改変問題の不祥事について謝罪と説明を行なうニュースについて書かれた Blog である。共通して出現した Blog は番組改変問題に関与した議院について触れているような、間接的にニュースに関わる Blog が見られた。対象語ではメディアのあり方を述べているような Blog が見られた。このことより、対象語による検索では番組改変問題というトピックを抽出できていると考えられる。対象語による検索は、そのニュースそのものについても抽出できると考えられるが、今回の実験では抽出が行なえなかった。これは、対象語によるキーワードの組み合わせに強い共起関係となるものが含まれにくいのに対し、頻出単語の共起関係にある単語が同時に上位に来る可能性が高いことが考えられる。今回の実験では、“橋本”と“会長”が強い共起関係にあるため、この組み合わせが用いられた検索では、頻出単語の方がトピックの特定としてふさわしいと考えられる。

- 検索結果数が極端に少ない結果となった。原因としては、ニュース中に使われる単語を利用した検索質問となるため、検索キーワードが Blog 中で使われるキーワードと異なっている場合や、そもそもニュース中に必要なキーワードが出現していないという可能性が考えられる。今回の実験では正解集合中によく出現する、“問題”、“改変”といったキーワードがニュース中に現れていない。そのため、そのような単語を含む質問を生

成できていないことが原因の一つだと考えられる。また“番組”といった重要であると考えられるキーワードがニュース中に含まれていたが、キーワード出現順が中ごろであり、重要度が高くないため抽出されなかった。ニュースとしては対象語として重要でないが、トピックとして重要と考えられるキーワードの抽出が行えないことは今後の課題である。

6. ま と め

本稿では、ニュース映像の評価と Blog の評価の統合方式として、ニュース映像のシーン順序により評価対象とその評価の抽出を行ない、その評価に基づく Blog 検索の提案を行なった。まず、対象語を定義し、プロトタイプの作成および評価実験を行なった。評価実験の結果として、対象語はニュースそのものではなくそのニュースのトピックを指し示すことができ、そのトピックについて評価を含むような Blog 検索に対して有効であることを確認した。

今後の課題としては、対象語で検索を行なう際、ニュースに出現するキーワードと Blog に出現するキーワードの違い、ニュースに含まれないがトピックとして重要なキーワードの抽出、評価を考慮した検索質問生成が挙げられる。またシーンの区切りの抽出がうまく行なえない場合には本方式が有効に働かないこと、シーン分割の自動化を実装し分割から検索まで行なえるシステムを構築すること、評価割合に関する評価実験方式の検討、現在別々に扱っている対象語と評価表現の対応付けの方式が挙げられる。さらに今回ニュース映像から評価の近い Blog を検索するという提案を行なったが、Blog からニュース映像を検索するという環境も考えられる。これらを検討することも今後の課題である。

謝 辞

本研究の一部は、平成 16 年度科研費基盤研究 (B)(2)「Web アーカイブと映像アーカイブを融合した次世代デジタル・ライブラリに関する研究」(課題番号: 16300028) によるものです。ここに記して謝意を表すものとします。

文 献

- [1] FNN-NEWS.COM
<http://www.fnn-news.com/>
- [2] WWW.NNN.COM
<http://www.nnn24.com/>
- [3] ANN NEWS
<http://www.tv-asahi.co.jp/ann/news/web/>
- [4] TBS NEWS i
<http://news.tbs.co.jp/>
- [5] Qiang Ma, Katsumi Tanaka, “WebTelop: Dynamic TV-content Augmentation by Using Web Pages”, Proc. of IEEE International Conference on Multimedia & Expo (ICME2003), Vol.2, pp.173-176, 2003.
- [6] Monika Henzinger, Bay-Wei Chang, Brian Milch, Sergey Brin, “Query-Free News Search”, Proc. of the Twelfth International World Wide Web Conference, 2003.
- [7] 若林正樹, 山田篤, 星野寛, 大瀬戸豪志, 上林弥彦, “評判情報の個人化によるサイト推薦システム”, DEWS2003 5-P-05
- [8] 奥村学, 南野朋之, 藤木稔明, 鈴木泰裕, “blog ページの自動収集と監視に基づくテキストマイニング”, 第 6 回セマンティックウェブとオントロジー研究会, SIG-SW&ONT-A401-01, 2004.

- [9] Tomoyuki NANNO, Toshiaki FUJIKI, Yasuhiro SUZUKI, Manabu OKUMURA, “Automatically Collecting, Monitoring, and Mining Japanese Weblogs”, Proc. of the Thirteenth International World Wide Web Conference, 2004.
- [10] Ravi Kumuar, Jasmine Novak, Prabhakar Raghavan, Andrew Tomkins, “On the Bursty Evolution of Blogspace”, Proc. of the Twelfth International World Wide Web Conference, 2003.
- [11] D.Gruhl, R.Guha, David Liben-Nowell, A.Tomkins, “Information Diffusion Through Blogspace”, Proc. of the Thirteenth International World Wide Web Conference, 2004.
- [12] 竹原幹人, 中島伸介, 角谷和俊, 田中克己, “Web 情報検索のための Blog 情報に基づくトラスト値の算出方式”, DEWS2004 I-2-02.
- [13] 中島伸介, 竹原幹人, 館村純一, 日野洋一郎, 原良憲, 田中克己, “blog 解析に基づく Web 情報検索の信頼性向上技術”, 人工知能学会第 6 回セマンティックウェブとオントロジー研究会, SIG-SWO-A401-05, 2004.
- [14] 中島伸介, 館村純一, 日野洋一郎, 原良憲, 田中克己, “リンク構造の時間特性に着目した Weblog 解析に基づくコンテンツ信頼性評価の検討”, DEWS2004 I-2-05.
- [15] 是津耕司, 日野洋一郎, 中島伸介, 門林理恵子, 呉受妍, 林正樹, 田中克己, “Weblog 情報を用いたコンテンツ・ブラウジング”, 人工知能学会第 6 回セマンティックウェブとオントロジー研究会, SIG-SWO-A401-04, 2004.
- [16] 藤村滋, 豊田正史, 喜連川優, “Web からの評判および評価表現抽出に関する一考察”, 情報処理学会研究報告 2004-DBS-134(II)-63, pp.461-468, 2004.
- [17] 藤村滋, 豊田正史, 喜連川優, “電子掲示板からの評価表現および評判情報の抽出”, 第 18 回人工知能学会全国大会, 3F1-03, 2004.
- [18] 藤村滋, 松村真宏, 岡崎直観, 石塚満, “電子掲示板上の評判情報に基づく意思決定支援”, 第 17 回人工知能学会全国大会, 2B1-05, 2003.
- [19] 鈴木泰裕, 高村大也, 奥村学, “weblog を対象とした評価表現抽出”, 人工知能学会第 6 回セマンティックウェブとオントロジー研究会, SIG-SWO-A401-02, 2004.
- [20] 立石健二, 石黒義英, 福島俊一, “インターネットからの評判情報検索”, 情報処理学会研究報告, NL144-11, pp.75-82, 2001
- [21] 形態素解析システム茶釜
<http://chasen.naist.jp/hiki/ChaSen/>
- [22] goo ブログ
<http://blog.goo.ne.jp/>