

# 視聴者の時刻同期コメントを用いた動画の特徴シーンの推定

末永 智彦<sup>†</sup> 早川 智一<sup>††</sup> 疋田 輝雄<sup>††</sup>

<sup>†</sup> 明治大学大学院理工学研究科基礎理工学専攻 〒214-8571 神奈川県川崎市多摩区東三田 1-1-1

<sup>††</sup> 明治大学理工学部 〒214-8571 神奈川県川崎市多摩区東三田 1-1-1

E-mail: †{suenaga,t\_haya,hikita}@cs.meiji.ac.jp

あらまし 本論文では、動画共有サービスに投稿された動画内容を効率的に把握するため、動画内の特徴的な盛り上がりを見せるシーンを発見し、それがどのようなシーンなのかを表す特徴語を推定する手法を提案する。実際に動画を見た視聴者の投稿したコメントは動画内容を反映していると仮定し、シーンの発見と特徴語の推定に用いる。例えばサッカーの試合の動画では、スーパープレーやゴールがあったシーンに「すげええ」や「うおおおお」といったコメントが集中する。提案手法では、ニコニコ動画に投稿された動画を対象に、機械学習で作成したネガポジ分類器によりコメントの感情を推定し、各感情のコメントが集中しているシーンを特徴的な盛り上がりのシーンとし、そのシーン内の特徴的なコメントを特徴語とした。提案手法の推定結果と実際に動画を見た評価者の回答との比較評価を行った。キーワード 動画共有サービス、自然言語処理、動画要約

## 1. はじめに

近年、動画共有サービスの普及と発展に伴い、ユーザが動画共有サービスの利用に使う時間が大きくなっている。また、動画を作成し Web へアップロードすることが容易となり、より多くの人が動画共有サービスに動画を投稿できるようになったことで、Web 上の動画の数も増加している。Youtube<sup>(注1)</sup>に1分間で投稿される動画の総再生時間数は、2017年時点で400時間[1]とされ、comScore<sup>(注2)</sup>によると、日本人1ユーザあたりの1ヶ月間の動画共有サービス利用分数は、2015年11月時点で2747分を記録した[2]。

これを受けて、ユーザが効率的に動画共有サービスを利用し、目的の動画をより短時間で視聴できる手法が必要とされていると我々は考える。

本研究では、ユーザの動画視聴時間を削減することによる効率化を目的とし、そのために動画の特徴シーンをを用いる。ここで、特徴シーンとは、動画内で特徴的な盛り上がりを見せる重要なシーンを意味する。目的の動画における特徴シーンの再生時間帯と、それがどのようなシーンなのかを表す言葉をユーザに提示することで、ユーザが動画の内容を効率的に把握できるようになり、視聴時間の削減につながると我々は考える。

しかし、特徴シーンの推定には、動画の特徴シーンがどこにあるか(動画がどこで盛り上がるか)は実際にユーザが動画を見てみなければわからない——という課題がある。

この課題を解決するため、動画に付随する情報を用いて自動的に動画の特徴シーンの再生時刻とそのシーンを表す特徴語とを推定する手法を提案する。動画のどの再生時刻が特徴的で、そこはどのようなシーンなのかをユーザが動画を視聴する前に提示することで、ユーザは視聴するシーンを取捨選択でき、視聴

時間の削減に繋がると我々は考える。

本研究では特に、動画に付随する情報として時刻同期コメント[3]を利用する。ここで、時刻同期コメントとは、動画の再生時刻に同期して投稿されたコメントのことを指す。時刻同期コメントを投稿できる点に着目し、日本の動画共有サービスであるニコニコ動画<sup>(注3)</sup>の動画を対象として、提案手法のプロトタイプの実装および実験と評価を行う。

時刻同期コメントは、「作品全体に対する通常のコメントに比べ、時刻同期コメントはその内容が質的に異なり、その瞬間ごとの感情をより多く表現する傾向がある」[3]と指摘されており、これを用いることで動画の特定のシーンの盛り上がりやその内容が推定できると考える。

特に、本研究ではコメントのネガポジ感情に注目する。例えばサッカーの試合を映した動画では、スーパープレーがあったシーンやゴールシーンがあった再生時間帯に対して、「すげええ」や「うおおおお」といったコメントが集中して投稿される。この場合、それらのシーンは多くのユーザが興奮した、つまりポジティブに盛り上がったシーンと言える。コメントのネガポジ感情推定は、サポートベクターマシン(SVM)を用いた教師付き学習によるコメントの感情分類により行った。

また、あるシーンの中で特徴的なコメントはそのシーンの内容を表していると考えられるため、それをそのシーンの特徴コメントとして利用する。例えば「すごいゴール」といったコメントが、動画内のあるシーンに集中して現れるなら、「すごいゴール」がそのシーンがどのような内容なのかを表していると言える。

我々の提案手法の特徴は次の3点である：

(1) SVMを用いた教師付き学習により、ニコニコ動画のコメントをネガティブ・ポジティブ・ニュートラルの3感情に分類した。

(注1) : <https://youtube.com>

(注2) : <https://www.comscore.com>

(注3) : <http://www.nicovideo.jp/>

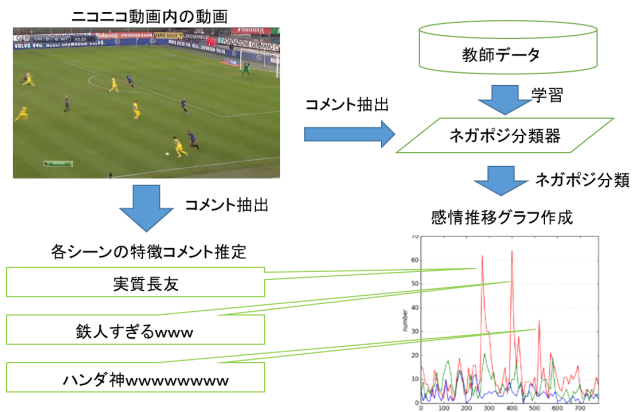


図1 提案手法の概要図

(2) 動画の特徴シーン推定を時刻同期コメントの時間的特性と(1)で分類した感情とを組み合わせを行った。

(3) 動画のあるシーン内のコメントから、そのシーンの内容を表す特徴コメントを推定した。

提案手法の概要図を図1に示す。提案手法のプロトタイプの実行結果をユーザが動画を視聴する前から得ることで、どのシーンを視聴すれば良いかがわかり、視聴時間の削減に繋がると考える。

今回作成した提案手法のプロトタイプは、次の2点をユーザに提供する：

(1) コメントの感情や集中度合いから動画の感情ごとの盛り上がり推移をグラフ化して特徴シーンを推定し、；

(2) 特徴シーン中のコメントから、そのシーンの内容を表す特徴コメントを推定する。

提案手法のプロトタイプを用いた評価実験では、特徴シーン推定の正答率は84%だった。評価実験は、ニコニコ動画に投稿されている動画に対しての実行結果と、実際に動画を視聴した人の回答とを比較して行った。

本稿の構成は以下の通りである。2章では関連研究について述べる。3章では本研究で対象とするサービスについて述べる。4章では提案手法について具体的に述べる。5章では評価実験について述べる。6章で本稿をまとめ、今後の展望について述べる。

## 2. 関連研究

本章では、動画に対するユーザのコメントを用いた研究を本研究に近いものとして挙げる。動画を対象とした研究は、動画の検索(2.1節)と動画分析(2.2節)とに分けて考えられるので、各節でそれぞれについて述べる。

### 2.1 動画検索

我々の研究は動画要約や内容の分析を行うものであるため、動画検索を行う研究とは異なるが、時刻同期コメントを用いた研究は、動画検索の分野でも多く行われている。

佃ら[4][5]は、ニコニコ動画の時刻同期コメントを用いることで、通常の動画のタイトルやタグに基づいたキーワード検索ではない、特定の人物の動画内での活躍度合いに注目した動画

検索手法[4]と、ユーザが自分好みの動画のランキングを生成できるシステム[5]といった、よりユーザの意図を反映した検索が行える手法を提案した。コメントを視聴者の反応として扱う点是我々と同じだが、これらの研究では動画全体の内容や特定の人物の活躍を知るためにコメントを用いており、我々は動画の特定のシーンの内容や盛り上がりの推移を知るために用いているという違いがある。

磯貝ら[6]や齋藤ら[7]は、視聴者からのコメントを元に動画の短時間なダイジェスト動画を作成し、それをユーザに見せることでその動画を視聴するかどうかの意思決定を助ける手法を提案した。これらの研究は実質的に動画要約を行っており、重要なシーンの推定に視聴者からの時刻同期コメントを用いるという我々と同じアプローチを取っている。しかし、盛り上がり度合いとして用いているのは単純なコメント数や、インターネット上で「笑い」を意味するスラングである「w」のついたコメント数であり、我々はコメントをネガティブ・ポジティブ・ニュートラルといった感情に分類して扱うため違いがある。

### 2.2 動画分析

時刻同期コメントを用いた動画要約の研究は、我々の研究に最も近いものとして挙げられる。山内ら[8]は、視聴者の観点の入れ替わりに基づいた特徴的シーン抽出を行った。彼らの研究では、動画内で視聴者のポジティブ・ネガティブ感情が入れ替わるタイミングにはそのきっかけとなる出来事があると考え、そこを特徴的なシーンだと仮定している。視聴者の感情推定には時刻同期コメントを用いており、その点是我々と同じである。だが、コメントの感情推定にはコメント内の単語と単語の感情極性の評価表現辞書とを用いており、SVMを用いて感情推定を行っている我々とは手法が異なる。我々も以前から本研究と同様の目的の研究[9]を行ってきた。これもSVMを用いたコメントの感情分類を行っているが、本研究では分類時のコメントに対する処理に変更を加えていること、新たにあるシーンの特徴コメントの推定を行っていることに違いがある。

動画に付随するコメントだけでなく、音声や画像といった動画自体の情報から動画要約を行う研究も多く行われている。中村ら[10]は、音楽を扱う動画に対して、音響特徴量とコメントの感情とを組み合わせたサビ部分検出を行うことでサムネイル画像を生成する手法を提案した。また松原ら[11]は、画像特徴量の変化に基づいて動画を複数の区間に分割し、それぞれの区間で最もコメントが盛り上がっているシーンを用いた動画の支持的ようなくサムネイルを生成する手法を提案した。これらの研究でも我々と同様にコメントをいくつかのクラスに分類して用いているが、辞書やあらかじめ決めた単語とのパターンマッチングで分類をしているため、SVMを用いている我々とは手法が異なる。

## 3. 対象とするサービス

### 3.1 ニコニコ動画

ニコニコ動画は、2018年1月時点で投稿動画数1,500万件を



図 2 特徴的なシーンの例

超える、ダウンゴ<sup>(注4)</sup>の運営する動画共有サービスである。ニコニコ動画は動画に対するコメント機能が他の動画共有サービスと異なり、動画の再生時刻に同期したコメントを投稿できる。Youtube などの一般的な動画共有サービスでは、ユーザのコメントは動画の枠外にリストで並べられ、ユーザ自身がコメント内に動画の再生時刻を記入するなどしなければ、それらは動画全体へのコメントとして見られる。これに対しニコニコ動画では、ユーザの全てのコメントが動画の特定の再生時刻に対して投稿される。ユーザが動画を視聴すると、あるコメントが投稿された再生時刻になると画面上にそのコメントが流れる。

このため、動画内の特徴的シーンにはコメントが集中する(図 2<sup>(注5)</sup>は、野球の試合で決勝点が入った瞬間のシーン)。多くの人がコメントを投稿したくなるような出来事があったシーンには、その出来事に関するコメントが大量に投稿される。

本研究では、この点に着目し、時刻同期コメントを動画の特定のシーンに対する視聴者の反応として扱う。具体的には、コメントの感情をそのシーンを見ている視聴者の感情とみなし、動画の感情推移を推定する。また、あるシーンに特徴的なコメントはそのシーンの内容を表していると仮定し、特徴コメントとする。

### 3.2 ニコニコデータセット

本研究では、動画に投稿されたコメントのデータセットとして、国立情報学研究所により提供されているニコニコデータセット<sup>(注6)</sup>を用いる。ニコニコデータセットは、ニコニコ動画のサービス開始当初から 2016 年 8 月 31 日までに投稿された約 1,400 万件のメタデータと、それらに対する合計約 35 億件のコメントデータを含むデータセットである。動画のメタデータについては、動画に付いたタグと動画のコメント数とを使用する。コメントデータについては、時刻同期コメントとして扱うため、コメント本文とコメントが投稿された動画上の再生時刻情報とを使用する。

本研究では、データセットの中でもコメント数が 1000 件以

上ある動画を対象とする。コメントの集中度合いから特徴シーンを推定するには、ある程度以上多くのコメントが投稿されている必要があると考えたためである。

## 4. 提案手法

提案手法では、動画内の時刻同期コメントを感情分類することによって得られた動画の感情推移をグラフ化したものと、各シーンにおける特徴コメントとをユーザに提示する。感情推移グラフにおいて山となっているシーンが動画の特徴シーンと言える。また、ここで特徴コメントとは、あるシーンの内容を最もよく表しているコメントのことを指すものとする。コメントの感情分類には SVM を用いる。コメントの素性をベクトル表現として SVM へ入力するために、形態素解析を行い、コメントを Bag-of-Words モデルに変換する。

提案手法の具体的な流れは、以下の通りである：

- (1) コメントに正規化処理を施す(4.1 節)。
- (2) コメントを Bag-of-Words モデルに変換する(4.2 節)。
- (3) SVM によるコメントの感情分類を行い、それぞれの感情のコメント数推移から特徴シーンを推定し、グラフ化する(4.3 節)。
- (4) 特徴シーン中のコメントから、シーンの内容をよく表す特徴コメントを推定する(4.4 節)。

### 4.1 コメントの処理

コメントに対する前処理として、正規化を行う。提案手法ではコメントを形態素解析して扱うが、ニコニコ動画に投稿されるコメントは基本的に短文であり単語の表記揺れも多く、そのまま扱うことが難しいためである。このような表記揺れが起きている同じ意味の単語をそのまま別々の単語として扱うと、特徴シーンの推定や特徴コメントの推定において各単語の数がまばらになり、結果が不正確になる可能性がある。

提案手法では、佃ら[5]の手法に倣い、以下の手順でコメントを正規化する。

- (1) スペース及び「!」と「?」以外の記号文字(「+」「#」など)を除去する。
- (2) 半角の片仮名と英数字、記号とを、全角に変換する。
- (3) 平仮名及び片仮名の小書き文字(「あ」「ア」「っ」「ッ」など)を大文字に変換する。
- (4) 英語の小文字を大文字に変換する。

以上(1)から(4)までの処理により、例えば「すげえええええ!!!」というコメントは「すげえええええ!!!」という文字列に正規化される。佃ら[5]の手法は Brody ら[15]の手法に倣い、同じ文字が 2 回以上繰り返して記述されている場合に繰り返し回数を 1 回にするという処理を行っていたが、我々はこれを行わない。4.2 節で述べるコメントの Bag-of-Words モデル化に関係するためである。

### 4.2 コメントの Bag-of-Words モデル化

コメントを感情分類のための SVM への入力や、特徴コメント推定に用いる TF-IDF 値の計算をするため、コメントを Bag-of-Words モデルに変換する。Bag-of-Words モデルとは、文章における各単語の出現回数のみを考え、単語の出現順など

(注4) : <http://dwango.co.jp/>

(注5) : <http://www.nicovideo.jp/watch/sm6535703>

(注6) : <https://www.nii.ac.jp/dsc/idr/nico/nico.html>

は考慮しない単純なモデルである。このモデルをコメントのベクトル表現とし、SVM への入力とする。

Bag-of-Words モデルへの変換のため、コメントの形態素解析を行う。形態素解析には、日本語形態素解析エンジンである MeCab<sup>(注7)</sup>を用い、システム辞書として mecab-ipadic-NEologd [12] [13] [14] を用いた。さらに、「w」、「8」、「!？」と、「ああ」のように母音を繰り返している文字列を感動詞として辞書に追加した。ニコニコ動画には多くのスラングが存在するが、特に感情を表しているながら多くの場面で記述される「w」（笑いを意味する）、「8」（拍手を意味する）、「!？」（驚きを意味する）も扱う必要があると考えた。母音の繰り返しについては、ニコニコ動画のコメントには、「すげええええ」や「おおおおおおお」と言ったコメントのように、強い感情を表すために母音を繰り返して記述される傾向があるためである。Brody ら [15] は Twitter<sup>(注8)</sup> 上のコメントには例えば「cool」が「coooooo!!!」になるように、単語の意味を強調するために同じ文字を繰り返して記述される傾向があることを報告している。

分解されたコメントの各品詞のうち、名詞と、形容詞、動詞、感動詞とを Bag-of-Words モデルに含める。

ここで、コメント中に同じ単語が複数回出現しても、ひとつのコメントにおけるその単語の出現回数は 1 とする。つまり、変換した結果のベクトルは、全体の語彙中の単語がそれぞれコメント中に出現しているかどうかを 0 か 1 で表す。

#### 4.3 特徴シーンの推定

SVM による線形分類器で、コメントの感情をネガティブ、ポジティブ、ニュートラルの 3 つの感情に分類する。SVM モデルの作成には、scikit-learn<sup>(注9)</sup> を標準的なパラメータで用いた。学習用の教師データとして、実際のコメントにそれぞれの感情のラベルをつけたものを人手で作成した。今回作成した提案手法のプロトタイプでは、教師データには 50 件の動画から無作為に 100 件ずつ抽出した 5000 件のコメントを用いた。

感情別にコメント数の時間推移を測り、コメント数の最も多いシーン 3 つを特徴シーンとする。複数の感情のコメント数が多いシーンについては、重複はせずにコメント数が多い感情の特徴シーンとする。コメント数は 10 秒単位で測る。また、動画の冒頭と終わりには投稿者に対しての挨拶コメントが多数投稿されることが多く、それらのコメントは動画の盛り上がりとは関係のないものであるため、動画の最初の 10 秒と最後の 10 秒は除外する。

また、各感情ごとに対象の動画におけるコメント数の時間推移をグラフ化する。グラフにおいて、各感情のコメント数が出ている山の部分はその動画の特徴シーンを表している。

#### 4.4 特徴コメントの推定

特徴コメントの推定には、コメント中の単語の TF-IDF 値を求めて使うというアプローチを取る。動画全体を 10 秒毎に区切り、それぞれのシーン中にあるコメント群をひとつの文書集

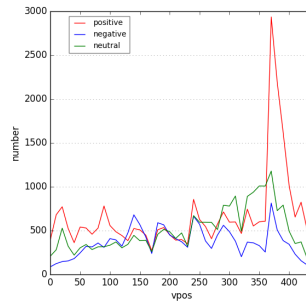


図 3 S1 の盛り上がり推移 (提案手法)

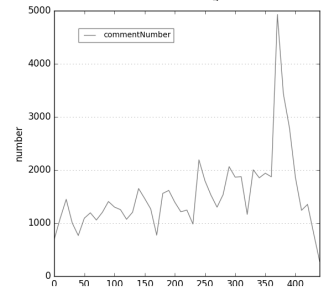


図 4 S1 の盛り上がり推移 (コメント数のみ)

合、動画全体のコメント群を全文書集合として単語の TF-IDF 値を計算する。例えば、動画内の特定のシーンにのみ「かっこいい」という単語が多く登場していた場合、そのシーン中の「かっこいい」という単語の TF-IDF 値は高くなる。

コメント中に含まれる単語の TF-IDF 値を合計をそのコメントのスコアとし、スコアがもっとも高いコメントをそのシーンの特徴コメントとする。この時、コメント本文があまり長い場合はそのシーンと無関係な内容について記述されていることが多いと考え、本文が 30 文字以下のコメントのみを対象とする。

## 5. 実験

### 5.1 実験方法

ニコニコ動画に投稿されている動画に対して提案手法のプロトタイプから得られた結果と、実際にその動画を視聴した 5 人の評価者<sup>(注10)</sup>からの回答とを比較して評価実験を行った。評価者はニコニコ動画のコメント表示機能をオフにした状態で動画を視聴した。特徴シーンの推定について、特徴コメントの推定について、同じ動画を対象としてそれぞれ評価を行った。

実験では 5 件の動画を対象とする。これらの動画はそれぞれ、ニコニコ動画内の「スポーツ」カテゴリの動画が 2 件、「ゲーム」カテゴリの動画が 2 件、「ニコニコ技術部」カテゴリの動画が 1 件である。

また、コメントの感情を考慮せずに単純にコメント数のみで動画の盛り上がりを推定する手法でも同様の実験を行い、提案手法と比較する。

### 5.2 特徴シーンの推定

実行結果の例として、「スポーツ」カテゴリの動画 2 件（それぞれ S1, S2 とする）について提案手法とコメント数のみの手法とで作成したグラフ、提案手法で推定した特徴シーンを図 3, 4, 5, 6 と表 1 に示す。図 3, 4 が、S1 に対してそれぞれ提案手法とコメント数のみの手法とで作成した盛り上がり推移のグラフで、図 5, 6 が、同様に S2 に対してのものである。各グラフは、縦軸が 10 秒あたりのコメント数、横軸が動画の再生時刻を表しており、提案手法で作成したグラフについては、赤線がポジティブなコメント、青線がネガティブなコメント、緑線がニュートラルなコメントの数を意味している。

例として挙げている動画の内容は、S1 が野球の国際大会の

(注7) : <http://taku910.github.io/mecab/>

(注8) : <https://twitter.com>

(注9) : <http://scikit-learn.org>

(注10) : 評価者の中に著者は含まれていない。

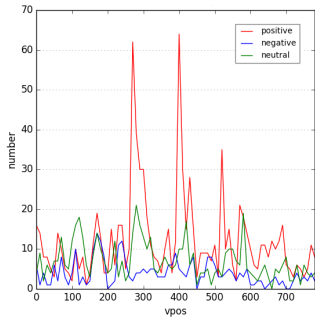


図5 S2の盛り上がり推移  
(提案手法)

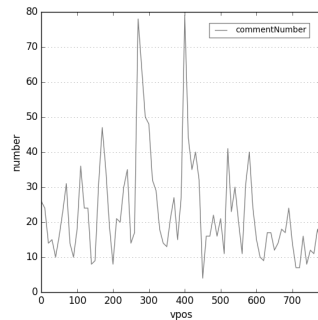


図6 S2の盛り上がり推移  
(コメント数のみ)

表1 推定された特徴シーン

S1		S2	
提案手法	コメント数のみ	提案手法	コメント数のみ
6分20秒	6分20秒	6分40秒	6分40秒
5分10秒	4分50秒	4分50秒	4分50秒
4分00秒	4分00秒	8分40秒	2分50秒

表2 S1の特徴コメント

再生時刻	特徴コメント
6分20秒	うおおおおおおおおおおおお
5分10秒	この場で今の球見逃せるのはすげーわ
4分00秒	とりあえず野球見ようぜ野球

表3 S2の特徴コメント

再生時刻	特徴コメント
6分40秒	鉄人すぎるwww
4分50秒	実質長友
8分40秒	半田神wwwwww

決勝戦で決勝点が入った打席の一部始終、S2が日本人選手も出場している海外サッカーの試合のハイライト映像である。

評価者は対象とする5つの動画を視聴して、もっとも重要だと思ったシーンを動画内の再生時刻(10秒単位)で回答した。それらの回答が、提案手法のプロトタイプが推定した3つの特徴シーンに含まれていれば正解とし、正解率を出した。

### 5.3 特徴コメントの推定

5.2節で提案手法により推定されたS1とS2との特徴シーンの特徴コメントの推定結果をそれぞれ表2、表3に示す。

全ての動画に対して表2、3のように推定されたコメントについて、「そのシーンの内容を表していると思うか」という質問を評価者に行い、評価者は「思う」、「思わない」の2通りで回答した。

## 6. 評価

### 6.1 特徴シーンの推定

提案手法で作成されたグラフとコメント数のみの手法で作成されたグラフと(図3と図4、図5と図6)を定性的に見比べると、見た目に大きな差は見られなかった。これは、提案手法の盛り上がり度合が、感情別とはいえ、コメント数のみの手法と同様にコメント数を元に測られているためだと考えられる。

表4 評価者の回答

評価者	動画 S1	動画 S2
A	6分20秒	9分30秒
B	6分20秒	4分30秒
C	6分20秒	4分50秒
D	6分20秒	4分30秒
E	6分10秒	6分30秒

表1より、結果に多少の差は見られた。

実験で示した例と同じ動画である、S1とS2とに対する評価者らの回答を例として表4に示す。

5件の動画に対する推定結果と評価者の回答とを比較すると、正答率は提案手法とコメント数のみの手法とでそれぞれ64%、56%だった。これは高い数値とは言えない。

しかし、表1と表4とを見ると、推定結果と回答とずれが10秒しかないものがある。実際に動画を見るとこれらは同じシーンを指していると考えられた。表4を見ると、動画S1に対してAからDまでは全員6分20秒が特徴シーンだと回答してEだけは6分10秒と回答しており、これも同様に同じシーンを指していた。

このように、時間の区切りや個人の感覚によって、同じシーンを指していてもずれが起こりうる。推定結果の前後10秒の回答を許容して正答率を測ったところ、それぞれ84%、76%に上昇した。提案手法の方が、8%高い正答率を得た。

動画ごとの正答率に注目すると、動画S1の100%に対して、動画S2は40%と低い。図3、図5を見ればわかるように、動画S1はグラフの山が動画全体で1箇所集中しているのに対し、動画S2は山が数箇所にある。また、動画S2の4分50秒部分や8分50秒部分を見ると、このコメントの盛り上がりは確かに動画の内容に関係のあるものだと考えられたが、普通に映像を視聴しているだけではそのシーンに起こっている出来事に気が付きにくいものだった。その出来事についてのコメントが表示されることによって気がついた他の視聴者がまたコメントを投稿していき盛り上がったシーンで、評価者はコメントを非表示にして動画を視聴したためにそのシーンの特徴に気が付かなかったものと考えられる。このように、動画を視聴する際にコメントを表示しているかしていないかによって、動画への印象が異なることが評価者の回答に影響していると考えられる。

### 6.2 特徴コメントの推定

推定された特徴コメントを「そのシーンの内容を表していると思う」と評価者が回答した割合は70.6%だった。

回答の中で特に「思う」という回答が少なかったのは、動画S1の4分00秒時点の特徴コメントである「とりあえず野球見ようぜ野球」(表2)であり、これは評価者AからEまでの全員が「思わない」と回答した。このシーンを見ると、コメント同士で言い争いが行われており、動画内容に関係があるコメントとは言えなかった。

逆に、動画S2の4分50秒時点と8分40秒時点と(表3)についてはどちらも5人全員が「思う」と回答した。これらは、

