

コメントの親子関係を利用した 動画共有サイトにおけるネットいじめコメントの検出

李 子怡† 川本 淳平†† フォン ヤオカイ†† 櫻井 幸一††

†九州大学大学院システム情報科学府 〒819-0395 福岡市西区元岡 744

††九州大学大学院システム情報科学研究院 〒819-0395 福岡市西区元岡 744

E-mail: †2IE15042N@s.kyushu-u.ac.jp, ††kawamoto@inf.kyushu-u.ac.jp fengyk@ait.kyushu-u.ac.jp,
†††sakurai@csce.kyushu-u.ac.jp

あらまし ネットいじめは、被害者に悪影響を与えることから、心理的に未熟な青少年の間で深刻な問題になっている。そして、近年ソーシャルネットワークサービス (SNS) の急速な発展に伴い、ネットいじめのインパクト面も広がりつつある。SNS における投稿からのネットいじめを検出するテキストベースの既存研究は、大抵研究者が索引語の集合を作り、コメントを単語単位で比較照合する。これらの手法は研究者の主観に影響され、各被害者のいじめ投稿の傾向に対応できないことが多い。また、動画投稿サイトは投稿したビデオを通じて物事を伝える情報基盤であり、コメントの対象はビデオの内容と投稿者と分かれている。そのため、本研究では、コメントとその返事の親子関係をスライド形式で捉え、ユーザー間のインタラクションを用いてネットいじめコメントを検知し、研究者の主観を控える。そして、コメントの対象がユーザーであるもののみをネットいじめコメントだとし、建設的な意見を排除出来る手法を提案する。

キーワード ソーシャルネットワークサービス, ネットいじめ, テキストマイニング

1. はじめに

近年、ソーシャルネットワークサービス (SNS) のユーザーは急速に増えている。ユーザーは親密またはプライベートな情報を気楽に共有している。このようなコミュニケーションは、悪意を持つ人に SNS を乱用する隙を与えた。投稿されたメッセージには、罵倒や失礼な内容を含む場合があり、更にはネットいじめにエスカレートするケースもある。ユーザーが安心して SNS を利用するためには、ネットいじめに該当するメッセージを自動検知する方法が必要である。

ネットいじめとは、インターネット上で故意に他人を侮辱、脅迫、困惑させる、苦しめる攻撃である。通常、成人は SNS に存在する危険を意識し、より安全なコミュニケーションをとることができる。それに対して、未成年者は脅威への認識力が低いため、身体的と精神的に大きな悪影響を受け、生死に関わる事件に至る可能性が高い。Ditch the Label 組織が 3023 名のイギリス学生を対象として行った 2015 年度いじめ状況サーベイ (ABS; Annual Bullying Survey)^(注1) では、いじめにあった結果、30% が自殺する考えを持ち、29% が自傷、27% が授業に欠席、14% が摂食障害を患い、12% が家出したなどの直接的被害が報告されている。インターネット以前の伝統ないじめの被害者は、ネットいじめに比べてより厳しいいじめ行為を受け普通の生活に影響を感じる事が多い。それに対して、ネットいじめの被害者は、物理的ないじめ行為を受けることはない

が、社会生活が困難になり不安や意気消沈を感じる傾向にあることが報告されている [13]。

ネットいじめは伝統的ないじめよりも固執である。伝統的ないじめは学校などでの接触を必要とするが、ネットいじめはそれらの環境を必要とせず被害者に一刻の猶予も与えない。特に SNS では、被害者はオンラインコミュニティ全体の前で傷つけられる [2]。また、「Once on the Internet, always on the Internet」と言われるように、一旦インターネットにアップロードされたものは、いかなる形式であっても永久に存在する。そのため、伝統ないじめと違う形式で、オンライン行為も反復性を持っていると言え、被害者に大きな悪影響を与える [10]。ABS によれば、調査に参加した学生のうち 43% もの学生がいじめを経験したと述べている。また、図 1 に報告にあったネットいじめの被害頻度を示す。図中 1 は経験したことがない、5 は度々、10 は高頻度でネットいじめの被害を受けていることを表している。62% ものいじめ被害者がネットいじめにさらされており、そのうちの 9% は頻繁にネットいじめされていることがわかる。このように、ネットいじめは無視できない影響と範囲を持っており、深刻な社会問題となっている。

Dadvar らによると、SNS におけるネットいじめ対策は主に人手に依存している。各サイトに担当者がおり、サイトに投稿された情報を監視しネットいじめに関係する投稿を手動で削除している [14]。しかし、大量の情報を処理する必要があり、手が及ばないことがある。このような状況を改善するために、計算機を用いた自動検知が期待されている。

既存のネットいじめに関係する投稿の自動検知手法は、主に辞書ベースの手法 [2] [3] と、ソーシャル情報を用いた手法 [4] [12]

(注1) : <http://www.ditchthelabel.org/the-annual-bullying-survey-2015-is-here/>

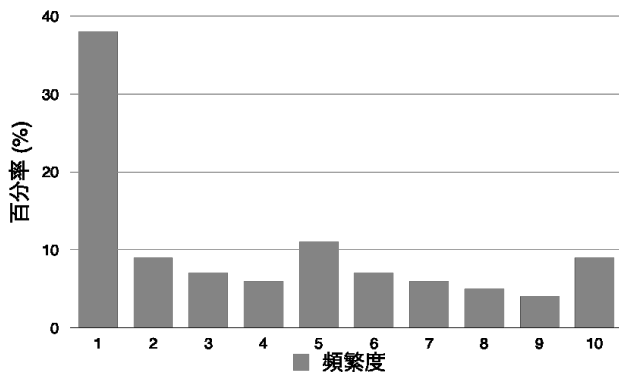


図1 ネットいじめの頻度。(Ditch the Label 調べ)

に分けられる。辞書ベースの手法は、簡単に実現できるが、いじめに関係する単語の選定や重み付け方法に辞書作成者である研究者の主観的な考えに影響されてしまう可能性がある。加えて、悪意ある単語を使わないネットいじめメッセージは検知することができない。特に英語は一つの単語でもたくさんの意味を持っており、インターネット上でさらに様々な意味が派生されているため、ブラックリスト辞書を用いた検出は難しい。ソーシャル情報、すなわちユーザ間の関係を表したソーシャルグラフを用いた手法やユーザ背景情報を用いた手法は、実際の投稿以外にもユーザ間のインタラクション履歴[4]やユーザの背景情報[12]など各種情報を必要とする。このような情報はユーザのプライバシーに関わるので、取得が比較的難しい。

本研究では、SNSのうち動画共有サイトを対象に、いじめ投稿を検出する方法を提案する。本研究で対象とする動画共有サイトでは、動画投稿者一人に対してコメントが集まる形式を考え、コメント投稿者は動画投稿者について背景知識を持っていると仮定する。Bastiaensensらの研究によると、ネットいじめにおいては、第三者はいじめ行為を目撃した際に被害者を助ける行動を取ることが多い[7]。実際、我々が動画共有サイトの一つであるYoutubeから取得したデータにおいても、第三者はいじめコメントに対して非難を向けるなど被害者救済行動が見て取れた。特にSNSでは、こうした第三者はいじめ被害者についての背景知識を持っていることが多く、どのような言葉が被害者にとっていじめに成り得るのかを考え、いじめ行為者へ非難を向ける、被害者を慰めるといった救済行動を行う。本研究では、こうしたネットいじめにおける人々の行動特徴を元に、単純な辞書式手法では発見できないが被害者にとっては傷つくような投稿を発見する。具体的には、コメントを単独ではなく、スレッドとして扱い親子関係を利用する。そして、いじめコメント投稿者と善意の第三者達とのインタラクションを用いてネットいじめコメントを発見する。また、SNSは投稿者を焦点とした短文や画像の共有する目的の他にも、物事を伝えるためにも使われている。そのため、コメントの対象はビデオの内容宛てと投稿者宛ての二種類が考えられる。そこで我々の提案手法では、コメントの対象を調べ、投稿者に対するいじめコメントを発見する。以降、本論文では英語のデータセットを用いて議論を行う。しかし、提案手法は言語に限らず使える方

法である。

本研究で提案された手法を用いて、少量のコメント情報と研究者の主観知識でネットいじめコメントを80%の真陽性率で検出することが出来た。ネットいじめコメントとして抽出されたものには、通常のいじめ話題ではなく、あるビデオ投稿者だけにとっていじめである話題であるコメントを含んでおり、既存のテキストベース手法では検知できないものを抽出できることが判明した。

本論文の構成は以下のとおりである。第2節では、SNSにおけるネットいじめ分析に関する先行研究を紹介する。第3節では、動画共有サイトにおけるコメントの扱い方を述べる。第4節では、本論文における提案手法を説明する。第5節では手法の実装を述べ、6節では実験の結果と考察について述べる最後に本稿をまとめる。

2. 既存研究

2.1 社会調査

ネットいじめ研究の初期段階では、ネットいじめの現象を理解するために、社会科学の専門家は、いじめの加害者と被害者の両方の心理的要因、人格と社会的関係に焦点を当ててきた。そのため、多くの大規模な調査を行い、いじめ現象のスコープを発見した。Bastiaensensらは、第三者のいじめに対する行動を調査した[7]。その結果、第三者はネットいじめを目撃した時に、被害者を助ける行動をとる可能性が高いことがわかった。ネットいじめの程度も第三者の行動と関連しており、いじめがより過酷であれば、援助行動を取る可能性も高くなる。一方、その第三者がいじめ加害者と友達であれば、被害者を助けるよりもいじめに加入する可能性が高いことがわかっている。

その他にも、ネットいじめを正しく定義すべく、伝統的ないじめとの違いに関する研究もある[10]。Langosらの分析によると、ネットいじめと伝統いじめの違いはIT経由であるかどうかのほか、反復性の影響も定義の一つである。同じ内容の反復は冗談、からかいと意図的な攻撃を区別する重要な基準であり、伝統的ないじめの定義には不可欠な要素である。一方、インターネット上では反復性へ捉え方が異なる。インターネットでは、どのような行為でもネットワーク上に無期限に存在し続けるという特性がある。すなわち、インターネット上の情報は簡単に転送でき誰もが閲覧可能である。したがって、たった一つのメールやメッセージであっても、その情報にアクセスされ閲覧されるたびに反復していると言えるのである。

Campbellらが9から19歳の学生3112名を対象として行ったサーベイでは、ネットいじめと伝統いじめの被害者の心理被害を調査している[13]。このサーベイでは、抑うつ不安ストレススケール(DASS, Depression Anxiety Stress Scales)を用いて、参加者の憂鬱、不安、ストレスのレベルを自己評価の形式で調査した。DASSは、大人の抑うつ、不安、ストレスの症状を測る尺度で、1項目0-3点の4段階で計算され、各3スケールの最低スコアが0で、尺度の最高スコアが42で、値の高さが問題と判断される。その結果を表1に示してある。いじめの経験者はいじめ経験なしの参加者と比べて、明らかにDASS

表 1 平均 DASS スコア

いじめの形式	憂鬱	不安	ストレス	総計
ネット	11.16	8.23	11.36	30.84
伝統	7.72	6.00	9.29	23.02
両方	14.62	11.73	15.35	41.70
経験なし	5.92	4.75	6.9	17.57

スコアが高い。そして、ネットいじめと伝統いじめ両方を経験している参加者は最も心理的な負担を抱えており、伝統いじめの被害者はいじめ経験者の中でより低い心理問題を持つことがわかる。

このように、ネットいじめに関する社会分析は数多く行われており、ネットいじめの悪影響と迅速な対処の必要性を示している。その一方で、ネットいじめ行為を自動的に検知する研究はいまだ少ない。

2.2 テキストコンテンツに対するネットいじめ検知

近年、コンピュータサイエンスの研究者を始めとし、ソーシャルネットワークサービスにおけるネットいじめの自動検知手法が提案されるようになった。これらは主に、SNS でのメッセージとコメントなどのテキストコンテンツを扱う手法である。インターネット上での会話にいじめ関連のキーワードが含まれているか調べることで、ネットいじめに関するメッセージを識別する Bag-of-words が基準線とされている [11]。

Yin らは感情分析と文脈の特徴を結合したモデルを提案し、Bag-of-words より良い精度が得られることを示した [11]。メッセージを tf-idf を用いたベクトル空間モデルで表し、感情情報を加えたものを特徴として用いることで、39.4% の精度、61.9% の再現率でネットハラスメントを検知した。Dinakar らは動画共有サイトの一つである Youtube に投稿されたコメントに含まれるネットいじめコメントを検知する方法を提案した [2]。彼らの手法では、まずコメントがいじめ関係する敏感トピック群に属するものか分類する。そして、さらにサポートベクターマシンを使ってどの敏感トピックに属するかを分類している。敏感トピックとは、性別、人種、知能や物理属性など人の簡単に変えられない特性である。彼らの手法は、ネットいじめであるコメントを検知できる精度は 66.7% である。Reynolds らは、コメントに対するラベル付けの中で、ネットいじめに関するメッセージは悪意ある単語を含む可能性が高いこと発見した [3]。そして、その情報を利用し、精度 81.7%、真陽性率 61.6% の自動分類器を提案している。単なる Bag-of-words ではなく、著者らは www.noswearing.com から得た 296 個の悪意ある単語をその悪意の度合いによりそれぞれ重み付けをした。各メッセージが含んでいる悪意ある単語の数及び密度を特徴とし、C4.5 決定木で分類している。

しかし、ネットいじめは社会的な現象であり、テキストで捉える情報だけでは不完全であることから、これらのテキストベースによる検出方法の精度は限られている。

2.3 ソーシャル情報を用いたネットいじめ検知

ここ何年か、テキストコンテンツのみを分析するのではなく、ネットいじめメッセージが交換された社会的な背景、ユー

ザーの背景や同時に投稿された画像を扱う手法も提案された。Huang らはソーシャルネットワークのグラフとしての特徴を分析することでネットいじめを検知する方法を提案した [4]。ノードとエッジ数、次数中心度、リンク数や k コアスコアなどを用いて、ネットいじめの加害者と被害者両方のソーシャル背景を特徴付けた。比較として、大文字、感嘆符や悪意ある単語の密度などのテキスト特徴も使用した。ユーザー間のソーシャルネットワーク構造の特徴とテキスト特徴両方を用いて分類した結果、ソーシャル情報を用いた手法はテキスト情報のみよる検知手法の検知率を改善できることが判明された。Huang らの手法における最も良い真陽性は 76.3% である。また、ネットいじめと検知されたソーシャルグラフを観察したのち、ユーザー間でのインタラクションすなわちリンクが比較的が多い場合は、ネットいじめであることも比較的に可能性が高いと述べられている。加害者と被害者の間では、予想以上にインタラクションが多いことも分かっている。Dadvar らはユーザーがネットいじめを受けた後の行為をネットいじめの自動検知手法の特徴として導入した [12]。ユーザーがネットいじめを受けたものとは異なるソーシャルネットワークサービスで行動を調査することで、ネットいじめの検知の精度を向上できることが示されている。同じ研究では、コメントの投稿者の性別もネットいじめであるかどうかの要因の一つであることを示した。女性は極端に冒涔な単語をそれほど使わず、関節的に否定的な単語使うことが多いとも示した。Kansara らは、画像とテキストの両方を分析し、ネットいじめを検知する手法を提案している [15]。画像とテキスト組み合わせは潜在的なネットいじめや脅迫の検知にとって良い情報の組み合わせであると述べられている。

以上の研究は、すべてコメントやメッセージは人を対象として、それぞれ単独な物であると仮定している。しかし、コメント間の関連は、ネットいじめメッセージが交換された際の情報を捉えられる。また、SNS は投稿者を焦点とした短文や画像の共有する目的の他にも、物事を伝える情報共有基盤でも使われている。物事に対しての意見と考えは人それぞれであり、誰もがその考えを伝える権利を持っている。そのため、ある物または事についてきつい事を言うのは建設的な意見とも考えられるので、全てのコメントやメッセージの対象を人だと仮定してネットいじめを検知するのは不完全である。

3. 動画投稿サイトにおけるコメント

本稿では、スレッド形式のコメント投稿を受け付ける動画投稿サイトを対象としている。本節では、そのデータ構造を定義する。

スレッドとはある特定の話題に関する投稿の集まりのことを言う。新たな話題を提供するためにコメントを投稿することを「スレッドを立てる」と言い、そのコメントをトップレベルコメントと呼ぶ。このトップレベルコメントに対する返信コメントや、返信コメントに対するさらなる返信コメントが連なりスレッドが形成される。そのため、トップレベルコメントと返信コメント間には親子関係があり、木構造の一種として扱うことができる。本節では、このコメント構造の形式的定義を与える。

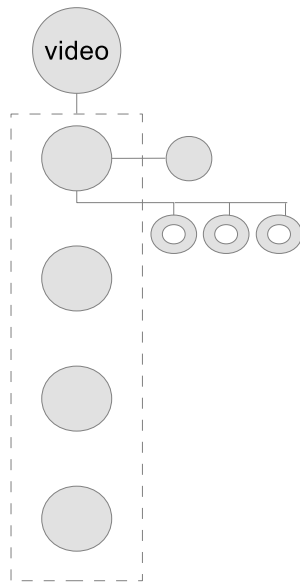


図2 スレッドを含むビデオ森の例 ○はトップコメント, ◎は返事投稿, は高評価数を示す

定義1: トップレベルコメント. コメント自身のテキスト情報とユーザーからの評価を含む半構造化されたマイクロブログのテキストである. C と示す.

定義2: 返事投稿. あるトップレベルコメントに対する返事. R と示す.

定義3: スレッド. トップレベルコメントとそのトップレベルコメントに対する返事投稿のセットである. スレッドは、ツリー構造として扱う.

$$T = (v, e)$$

と示す. そのうち,

$$v = \{c \cup \mathbb{R} \mid c \in \mathbb{C}, \mathbb{R} = c \text{ への返事集合}\}$$

$$e = \{\{a, b\} \mid \text{コメント } a \text{ と } b \text{ が関連}\}$$

定義4: データセット. スレッドからできる森と表せる.

$$F = (\mathbb{V}, \mathbb{E})$$

と示す.

例えば, 図2に示されているデータセットには, 四つのトップレベルがあり, すなわち, 四つのツリーからできる森である. 各ツリーには, トップレベルコメントがツリーの根あり, それに対する返事がノードである.

ツリー構造は, ノードにあるテキストデータを示すことができるほか, 重要な投稿間の関連もあらわにすることができる.

4. 親子関係を用いたコメント解析

本稿では, 動画共有サイトにおいてコメントの対象分岐を解決するために, 投稿コメントの対象を人と物に分けて扱い, またコメントの親子関係を利用したネットいじめ検知手法を提案

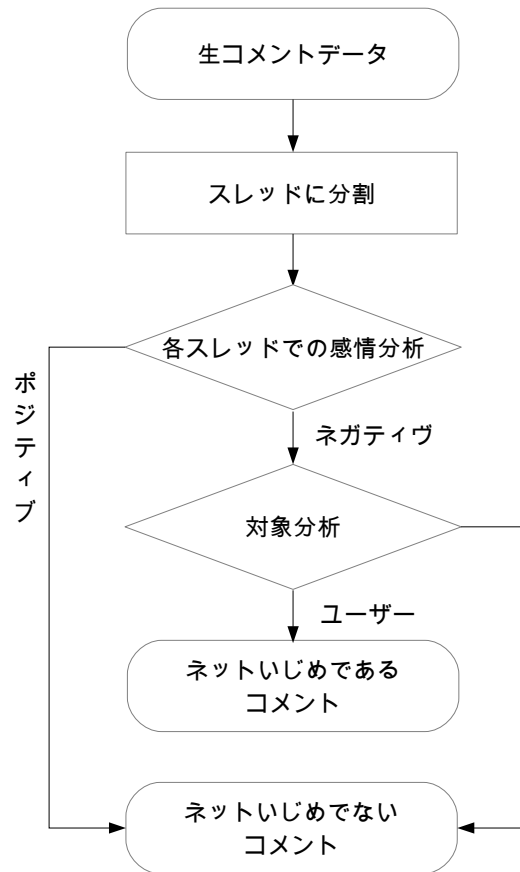


図3 提案手法の手順

する. 投稿コメントの対象を区別することで, 物事を対象とした建設的な意見をネットいじめであると誤検知することを防ぎ, コメントの親子関係を考慮することでネットいじめメッセージが投稿された際のユーザー間のインタラクションを捉える.

提案手法の手順を図3に示す. 動画投稿サイトから取得したコメントをスレッド単位に分割する. そして, 各スレッドごとに子コメントに対する感情分析を行う. 得られた結果を元に親コメントに対するスコアを求める. 最後に, 親コメントの対象を調べ, 対象がユーザーであるコメントをネットいじめコメントとして抽出する.

子スレッド投稿されたコメントは, 親コメントに対するポジティブまたはネガティブな反応を含んでいることがあり, ネットいじめにおいては, 善意の第三者はネットいじめに関係するコメントに対しては否定的な態度をとる [7] ことから, 子コメントが多くネガティブな反応を含んでいる場合, 親コメントはネットいじめに関する内容を含んでいる可能性が高いと判断できる. 特に, 動画の視聴者の一部は動画投稿者についてよく知っていると仮定できる. すなわち, その投稿者の動画を購読しているような視聴者は, 投稿者の好みやどのような言葉に傷つくのかといった背景知識を共有していると仮定する. 感情分析では, コメントをスレッド形式で捉え, コメント投稿者間の反応を解析することで視聴者からの反応を調べる.

コメントツリー

$$T = (v, e)$$

$$v = \{c \cup \mathbb{R} \mid c \in \mathbb{C}, \mathbb{R} = c \text{ への返事集合}\}$$

$$e = \{\{a, b\} \mid \text{コメント } a \text{ と } b \text{ が関連}\}$$

があるとする。R に対して感情分析を行い、返事の全体がネガティブ評価であるか、ポジティブ評価であるかを得る。その評価に基づき、ツリーの根である c を評価する。

返事文から得られる感情のほか、YouTube にある各トップレベルコメントについている高評価ボタンが押された回数も使用した。高評価ボタンは、図 4 に示されているように、コメント欄の下方に位置しており、コメントが有用、または同感を感じた時に使われる。図 4 から見ると低評価ボタンもあるが、Google社が YouTube を買い取った後は使われず、形式だけ残っている。高評価ボタンを押すことはトップレベルコメントを賛同しているため、ボタンが押された回数もポジティブな物として扱い、総合評価に足している。

対象分析では、コメントの対象を調べ、投稿者に対する悪意あるコメントのみを抽出する。物事を伝える情報基盤として使われているソーシャルネットワークサービスでは、コメントの対象はビデオの内容と投稿者とで分かれている。感情分析で得た低評価のトップレベルコメントを物事に対する建設的な意見であるか、またはユーザーに対するネットいじめである。すなわち、動画で語られている物事に対する否定的なコメントや一見いじめに見えるコメントは、それらの物事に対する意見であり、個人攻撃とは限らない。本稿では、動画投稿者へのいじめコメントの発見が目的である。そのため、対象分析を行い、投稿者宛てのコメントのみを抽出している。

5. 実装と評価実験

5.1 コーパス

本研究で用いたデータセットはソーシャルネットワークサービスの www.youtube.com から抽出したものである。YouTube とは Google社が運営する、世界で最も利用者が多い動画共有サービスである。YouTube という SNS をデータ源として選んだ理由は以下が挙げられる。

(1) YouTube 利用者の人口分布は一般インターネット利用者の分布と一致している [1]。そのため、得られた結果はこのサイトに限らず、インターネット上におけるユニバーサルなものである可能性が高い。

(2) ネットいじめであるコメントが存在する情報源である。消費者が内容を生成していくメディア消費者生成メディアであることから、ビデオ投稿者個人に対するネットいじめ行為や嫌がらせコメントを投げ出すプラットフォームとして乱用されるケースもある。YouTube は動画所有者に動画からコメントを削除する権利を与えますが、動画投稿者が実際にコメントを見た時点で被害を受けていると言えるため、ネットいじめ行為が緩和できるとは言えない。特に、論争のトピックに関するビデオは、多くの場合、不快と失礼なコメントがある。

(3) コメントの対象はユーザーに限らず様々である。投稿

された動画のトピックは幅広く、投稿者の生活動画はもちろんのこと、オンライン講義など物事を伝える情報基盤としても使われている。また、基本的には実名サイトではないので、人間関係は比較的薄く、ユーザー間の直接的なメッセージのやり取りは少ない。

(4) コメント量が豊富である。YouTube は 2013 年末から、スレッド形式のコメント欄を導入しており、10 億ものユーザーがいるため、十分なデータが得られる。

著者らの知る限りでは、ネットいじめ検知に使えるオープンな YouTube コメントのデータセットはない。そこで、独自に YouTube からコメントデータを取得した。本研究は、投稿者一人一人異なる不快な言説も含めていじめ関連コメントを発見することが目的である。そのため、セクシュアリティ、人種と文化と知性などの一般的に議論を招くものをトピックとした動画のほか、美容、ファッションに関する動画からもコメントを取得した。なお、動画投稿者が特定できないパブリックアカウントからの動画は排除した。また、YouTube 自体は多言語で構成されており、投稿動画の言語もそれぞれだが、本研究では英語でのコメントのみを取得した。

次に、得られた動画から、トップレベルコメント、コメントに対する返事と押された高評価ボタンの回数を取得した。本稿で使用するデータセットは、58 本の動画からなり、合計で 41415 件のトップレベルコメントを採集した。そのうち、2036 件のトップレベルコメントは子コメントを持ち、全体の約 5% を占めている。子コメントは全部で 8313 件であった。

本研究で使用した情報は、コメントのテキスト文とその高評価数である。YouTube を含む動画共有サイトからは一般的に、その他のデータも収集することができる。本研究では、最少限度のデータを使用し、それらの関係を考慮することでネットいじめコメントを分析できるかを知るべく、上記の情報のみを用いることにした。

スレッド形式のコメントの例として、YouTube の二つのトップレベルコメントとそれに対する返事を図 4 に示す。

図 4 に示す例では、動画視聴者 A と E がそれぞれスレッドを立て、そのコメントに対する返事コメントが投稿されている。もし、その返事が直接トップレベルコメントの投稿者に対するものでなければ、「+」記号の後にそのコメントの対象ユーザーを示している。しかし、現在の YouTube ではセカンドレベル以降のコメントに対する返事の取集を制限しているため、本研究ではスレッドの中の返事は全てトップレベルコメントに対するものとする。また、コメント投稿者は別のスレッドを立てて、他のユーザーのコメントについて評価する事もできる。しかし、比較的少ない行為であるため、今回は独立したコメントとして扱う。我々が調べたところ、実際には無関係のコメント投稿、例えば、スパムや絵文字などの本研究にとって無意味な投稿が存在するが、極性分析には影響が少ないことが分かっている。

5.2 感情分析

OpinionFinder (OF) は一般に入手可能なソフトウェアであり、文脈レベルの主観性を識別し、感情の二面性であるポジティブ・ネガティブを判定することで感情分析を行う [5]。この



図 4 スレッド例

ツールによる広範なツイート集合の感情分析は既に行われており [6]、その研究では、OF により得られた時系列の気分データが、その調査期間において、株価を予測の正確性を向上することを示した。OF の感情分析は、OF の「辞書」により、ツイート内容を解析し、ツイートの「ポジティブ/ネガティブ」を決定する。そこで、我々は過去の研究 [8] にて実績のある OF 辞書を今回の検証に採用した。採用した OF 辞書は 2718 個のポジティブワードと 4912 個のネガティブワードを持つ。この辞書では、さらに単語を「強気」もしくは「弱気」のポジティブとネガティブに分けてあるが、ネットいじめは通常強い感情と関連付けられているので、本実験では強気な単語のみ評価に入れた。そして、それぞれのツイートにおいてネガティブかポジティブのどちらの感情が支配的かを求める。強気なネガティブを表現する言葉数と強気なポジティブを表現する言葉数を求めた。

OF での辞書には、攻撃的な単語は含まれていない。先行研究のラベリングプロジェクトで、ネットいじめにレベル付けされやすいメッセージは攻撃的な単語を含む可能性が高いことが示された [3]。このようなトップレベルコメントを反論するときも、悪意ある単語を使う可能性があるため、我々はウェブサイト www.noswearing.com に掲載してある攻撃的な言葉のリストを更に拡張し、350 個の悪意ある単語とそれらの変形を強気なネガティブ単語として扱った。

以上 OF から求められたポジティブ数に高評価ボタンが押された回数を足し、ポジティブ総数とネガティブの総数が得られる。そして、多い方が支配的な感情とすることにより、トップレベルコメントに対する返事を感情で二分化し、返事の感情を計算した。

5.3 対象分析

エンティティとエンティティレベルの感情の抽出のために我々は AlchemyAPI を使用した [9]。AlchemyAPI は、ディープ・ラーニング、自然言語処理技術及び機械学習アルゴリズムを利用して、リアルタイムのテキスト分析、コンピューター・ビジョンを API として提供している。AlchemyAPI を使用するために API 鍵を必要とするが、それは無料で入手できる。ただし、1 日あたり 1000 回の API 呼び出し回数制限がある。AlchemyAPI を使用している開発者は 36 カ国、4 万人以上いると言われている。

本実験で使用するのは、Alchemy Language である。このサービスは、キーワード、名前付きエンティティ、およびエンティティレベルの感情を抽出することができる。AlchemyAPI のエンティティ抽出では、テキスト内の人々、企業、団体、都市、地理的特徴、および他の型指定されたエンティティを識別することができる。さらに、28 種類のエンティティのほか、100 種類以上のサブカテゴリもテキスト文から抽出できる。統計的なアルゴリズムと巨大なデータベースを組み合わせ、色々なオブジェクト、個人、および場所を記述する。また、政治家やアスリートとして個人を識別すれば、詳細なオントロジーマッピングも提供できる。以上の機能は、本研究にとって有用である。

対象分析として、各トップレベルコメントに以下の操作を行った：

- (1) 重要でない文字列の除去：「lollllll」、「HAHAHAHAHA」と「Wow」などの文字列はデータセットから除去した。
- (2) 略記を書き直す：オンラインでよく使われている略記を元書き直す。例えば、「u」を「you」に直し、以後の処理を容易にする。
- (3) 人称代名詞の書き換え：AlchemyAPI ではエンティティが代名詞で表されている場合は抽出されないため、人称代名詞をビデオ投稿者のユーザー名に書き換えた。ある特定のトップレベルコメントの人称代名詞を書き換えるかどうか、複数の人称代名詞が存在する場合はどれを書き換えるかを定めるルールはとても難しい。そのため、本研究では、トップレベルコメントの第 1 句で一番最初に出現する人称代名詞をユーザー名に書き換えた。「you」、「he」と「dude」などを「Alice」と「Bob」のように入れ替え、AlchemyAPI で処理できるようにする。
- (4) AlchemyAPI でのエンティティ分析。
- (5) 直接ビデオ投稿者に対するネットいじめコメントを抽出対象とする。有名人、ある製品や地域などをエンティティとしたトップレベルコメントは取り除き、エンティティがビデオ投稿者であるもののみをネットいじめコメントとして登録する。

5.4 評価方法

ネットいじめだと感じるかどうかは被害者によりますが、コメントの投稿者は他人を傷つける意図を持っていない [10]。したがって、ビデオ投稿者へのポジティブコメントはネットいじめではないと想定できる。そのため、本研究では、AlchemyAPI のエンティティレベルの感情分析で得たスコアを

基準として、抽出されたネットいじめコメントの真偽性を評価する。エンティティレベルの感情スコアがプラスである場合は、エンティティにポジティブな感情を伝えているので、偽ネットいじめであるとし、その他を真ネットいじめコメントとする。

6. 考察

実験で得た真陽性率を表 2 に示す。伝統的な精度という評価基準は、両クラスが特に不平衡である時や誤分類のコストが両クラスで非常に大きな違いがある時は分類器をうまく評価できない。例えば、一番簡単な分類器 ZeroR を使い、全てのコメントを多数クラスと想定する場合でも、ネットいじめコメントの数が全体のコメント数と比べて非常に少ないため、100%に近い分類性能が得られる。そのため、本研究では、提案手法の性能を評価するため、少数クラスの真陽性率を使用する。先行研究でも、少数クラス（重要な）クラスに注目し、性能を図っていることが多い [4]。

表 2 TP

	スレッド分析のみ	対象分析
TP	36 %	80 %

対象分析は、スレッドごとの感情分析で得たネットいじめコメントの真陽性率を向上したことが表 2 で観測できる。スレッドごとの感情分析のみで抽出されたコメントリストでは、ビデオで討論された話題に関するものや人が対象であるコメントが多く含まれている。例として、古代残酷な統治者についてのビデオでは、「Caligula was not insane.. he was just an asshole to his consuls and the senate.」というコメントが残された。ビデオ内容には反論し、悪意ある単語も含まれているが、ネットいじめコメントとは言えない。

抽出したリストの中には、悪意ある単語を含まないネットいじめを含む。表 3 に示されているトップレベルコメントを例とする（ビデオ投稿者のユーザー名を控えるため、ここでは Alice とする）。このコメントの総合評価は -18 である。そのうち、返事 1 は -2 であり、「bye」を含む返事 2 から 4 はそれぞれ -5 との評価を得た。四つの返事の中では、返事 1 が特に面白く、評価としては -3 である。返事 1 では、一見悪意を持たないトップレベルコメントの裏側の意味を展開し、投稿者をサポートしつつ、トップレベルコメントの投稿者を批判している。返事 1 の前 3 句は全体的にニュートラルであり、最後の一句「Your ignorant, and rude」がネガティブと判断したキーである。

伝統的な、対面のいじめでは、同じ内容の反復は冗談、からかいと意図的な攻撃を区別する重要な基準であり、伝統的ないじめの定義には不可欠な要素である [10]。一方、サイバースペースの特質は反復性への理解を変えた。サイバースペースにおける一つの行為でも、ネットワークに無期限に残される特性を持っていながら、簡単に転送でき、誰もが見られる。したがって、一つのメール、メッセージ、ポストであっても、アクセスや見られた度に反復していると言える。表??にあるトップレベ

表 3 単独なネットいじめコメント例

トップレベルコメント	you get paid for doing this, so im not watching anymore. you are paid and we are tricked into buying products,	-18
返事	Yeah so, but her opinions are true and honest. We don't need people like you on Alice's channel. Alice spends a lot of money and Bob too on things. They need to make money to survive, but Alice isn't a liar. You aren't tricked into buying these products that are being sponsored. Your ignorant, and rude	-3
	bye	-5
	bye	-5
	lol bye	-5

ルコメントの内容はこのコメントにしかないが、ビデオ投稿者がコメントをチェックする度にこのようなネガティブにさらされているのである。

単独に出現するコメントのほか、一人のビデオ投稿者に対して、同じトピックのトップレベルコメントもリストに存在する。表 4 に示されているトップレベルコメントは、それぞれ違うユーザーから投稿されたものであるが、どれも同じ話題についてのものである。トップレベルコメント 1 の評価が特に低い原因は、返事の中に悪意ある単語が複数含まれているからである。もし、悪意ある単語を除けば、これらのトップレベルコメントの評価はどれもそれほど低いわけでもない。しかし、ビデオ投稿者のチャンネルではこのようなコメント内容に対する不満を示すビデオが実際にアップロードされている。たとえニュートラルに近いコメントであっても、出現する回数の上昇により、不快感を感じる可能性も高まる。

表 4 同じ話題であるネットいじめコメント例

トップレベルコメント	返事数	評価
dude put weight on	100	-68
it's been a while since i watched wayne channel.. but is he all right?!?! is he ill?!?!?!?! i hope hes is fine because he looks horrible ... anyway thanks for this video Bob...	4	-5
You lost weight, are u doing anything?	1	-2
what happened to u ? r u sick or u just loose ur weight ?	1	-1

その中で、トップレベルコメント 1 は平均返事数が 4 件の中で、100 件の返事を得ている。返事をぐたいに観察すると当コメントを投稿したユーザーがその他のユーザーと言いつていることがわかる。同じ原因で悪意ある単語の出現頻度が高くなっている。「痩せた」という内容は、一人一人違う考えを持つので、場合によれば、特に女性にとっては嬉しいことである。そのため、意見が分かれており、言い合いするようになった。また、今回のデータセットにはタイムスタンプがないので、ウェブ上で見てみれば時間に連れ、「痩せた」という内容のコメントに対する返事の感情が変わっていることがわかる。ビデオ投稿者は近日実際に痩せているので、彼が不満を示すビデオを

アップロードする前に、実際に多くの人がこのトップレベルコメントに賛同していることがわかる。「痩せた」との内容のコメントに不満を示したビデオのアップロード後、多くのビデオ投稿者の購読者がこのような内容のコメントに否定する返事をするようになった。

7. まとめ

本研究では、近年 SNS 上で問題視されているネットいじめの自動検知手法として、コメントの対象とコメント間の親子関係に着目した。トップレベルコメントとその返事との親子関係を考慮したのち、辞書ベースでは検知できない各被害者に対して特定話題のネットいじめコメントを抽出することが出来た。また、コメントの対象をユーザーのみのものに絞ることで、ある話題への建設的な意見を排除でき、80% の真陽性率を得た。

本稿で提案された手法は、コメントのテキスト以外の情報と研究者の主観への要求が少ないのが利点である。一方、不足部分もいくつかある。コメントの返事を基準として、コメント自身を評価しているので、返事を持たないネットいじめコメントは無視されてしまう。また、返事の中でのネットいじめコメントは検知できない。これらの不足点を解決するのが今後の課題である。本研究で使われたのはコメントのテキスト文とその高評価数であるが、明らかに、YouTube サイトからさらなる豊富なデータを収集することができる。データフィールドを拡大し、プロフィール情報などを用いてツリー構造を改善し、より細かなコメントの分析に対応させる。

8. 謝辞

本研究の一部は、九州大学 P & P つばさプロジェクト採択研究「IT 技術を用いた大学生対象のメンタルヘルス e-ラーニングシステムの構築」支援を受けている。

文 献

- [1] Cha, Meeyoung, et al. "I tube, you tube, everybody tubes: analyzing the world's largest user generated content video system." Proceedings of the 7th ACM SIGCOMM conference on Internet measurement. ACM, 2007.
- [2] Dinakar, Karthik, Roi Reichart, and Henry Lieberman. "Modeling the detection of Textual Cyberbullying." The Social Mobile Web. 2011.
- [3] Reynolds, Kelly, April Kontostathis, and Lynne Edwards. "Using machine learning to detect cyberbullying." Machine Learning and Applications and Workshops (ICMLA), 2011 10th International Conference on. Vol. 2. IEEE, 2011.
- [4] Huang, Qianjia, Vivek Kumar Singh, and Pradeep Kumar Atrey. "Cyber Bullying Detection Using Social and Textual Analysis." Proceedings of the 3rd International Workshop on Socially-Aware Multimedia. ACM, 2014.
- [5] Wilson, Theresa, et al. "OpinionFinder: A system for subjectivity analysis." Proceedings of hlt/emnlp on interactive demonstrations. Association for Computational Linguistics, 2005.
- [6] Bollen, Johan, Huina Mao, and Xiaojun Zeng. "Twitter mood predicts the stock market." Journal of Computational Science 2.1 (2011): 1-8.
- [7] Bastiaensens, Sara, et al. "Cyberbullying on social network sites. An experimental study into bystanders' behavioural

intentions to help the victim or reinforce the bully." Computers in Human Behavior 31 (2014): 259-271.

- [8] Wilson, Theresa, Janyce Wiebe, and Paul Hoffmann. "Recognizing contextual polarity in phrase-level sentiment analysis." Proceedings of the conference on human language technology and empirical methods in natural language processing. Association for Computational Linguistics, 2005.
- [9] "Alchemy api," . [Online]. Available: www.alchemyapi.com
- [10] Langos, Colette. "Cyberbullying: The challenge to define." Cyberpsychology, Behavior, and Social Networking 15.6 (2012): 285-289.
- [11] Yin, Dawei, et al. "Detection of harassment on web 2.0." Proceedings of the Content Analysis in the WEB 2 (2009): 1-7.
- [12] Dadvar, Maral, et al. "Improved cyberbullying detection using gender information." (2012).
- [13] Campbell, Marilyn, et al. "Victims' perceptions of traditional and cyberbullying, and the psychosocial correlates of their victimisation." Emotional and Behavioural Difficulties 17.3-4 (2012): 389-401.
- [14] Dadvar, Maral, et al. "Improving cyberbullying detection with user context." Advances in Information Retrieval. Springer Berlin Heidelberg, 2013. 693-696.
- [15] Kansara, Krishna B., and Narendra M. Shekoker. "A Framework for Cyberbullying Detection in Social Network." International Journal of Current Engineering and Technology 5.1 (2015).