

## 機械学習を用いた自然言語処理による商品レビューの評価<sup>†</sup>

市川 知春<sup>\*1</sup>・武田 和大<sup>\*2</sup>・原 崇<sup>\*2</sup>

### Evaluation of Product Reviews by Natural Language Processing Using Machine Learning

Tomoharu Ichikawa<sup>\*1</sup>, Kazuhiro Takeda<sup>\*2</sup> and Takashi Hara<sup>\*2</sup>

**Abstract** The review function of the internet shopping is the function which the purchaser can write impression and evaluation of the goods. It is possible that the purchase applicant reads the product review, and that it obtains information of the product beforehand. However, some product reviews are not helpful because all buyers do not write reviews carefully. In this study, we analyze and evaluate product reviews by natural language processing using machine learning, and construct a system that rearranges products in order of reference. Learning is carried out using logistic regression so that the words which appears in the content of the review many times is made to be the learning feature, and the appropriate weight in which the learning feature affects the content of the product review is calculated. For the evaluation of the system, the QE method which is the evaluation method using the quick sort is proposed. By using the QE method, it was possible to evaluate whether the rearrangement was right. As a result of the experiment, it was confirmed that the proposed system could rearrange the product reviews with high accuracy.

**Key words** Machine Learning, Natural Language Processing, Logistic Regression, Quick Sort

#### 1. はじめに

近年、インターネットを通じて商品やサービスを購入できるインターネットショッピングの需要が高まっている。インターネットショッピングの多くには、購入者が商品の感想や評価を記入できるレビュー機能が備わっている。例えば、世界最大級のインターネット通販サイト Amazon では、購入者が商品ごとに文章による商品レビュー（以下、レビュー）と5段階の数値による評価を投稿することができる。購入希望者はその商品レビューを読むことで、事前に商品の情報を得ることができる。さらに、Amazon の場合は「役に立った」ボタンが備えられており、そのレビューが参考になる情報であるのかどうかを判断することができる。しかし、商品の中には、数千件のレビューの中で、半数程度が「役に立った」という評価がされてい

ない場合があった。つまり、数多くのレビューがあっても、その中には参考になるレビューもあれば、参考にならないレビューも存在しているということである。例えば「届くのが楽しみだ」ということで、その商品を高く評価している場合や、「体に合わないサイズを買ってしまった」ということで、その商品を低く評価している場合は、その商品の情報ではないため、参考にならないと言える。購入者が投稿したレビューが参考になるかならないかは、Amazon の「役に立った」ボタンのようにショッピングサイト自体の機能で判断できるものもある。しかし、このようなレビューを評価する機能を備えていないショッピングサイトも存在する。この場合は、購入希望者は数多くあるレビューを読んで、その中から購入の参考になる情報が書いてあるレビューを探す必要がある。

このような背景から、自動的に参考になる順序にレビューを並び替える技術が求められるが、レビューの内容は自然言語で記述されることが普通であるため、その解析は容易ではなく、有効な方法は提案されていない。

そこで、本研究ではショッピングサイトのレビュー

<sup>\*1</sup> 鹿児島工業高等専門学校 専攻科 電気情報システム工学専攻  
Advanced Electronic and Information Systems Engineering,  
National Institute of Technology, Kagoshima College

<sup>\*2</sup> 鹿児島工業高等専門学校 情報工学科  
Department of Information Engineering, National Institute of  
Technology, Kagoshima College

<sup>†</sup> 2021 年 1 月 25 日受付・2021 年 7 月 21 日再受付

について、機械学習を用いた自然言語処理にて解析、評価を行い、参考になる順番に並び替えを行うシステムを構築することを目的とする。具体的には、ショッピングサイトからレビューを取得する機能と、取得したレビューを学習して参考になる順番に並び替えを行う機能を開発し、システムを構築する。さらに、並び替えた順番が正しいかどうか評価するために、新たな評価法となるQE法を提案する。最後に、提案するQE法で、開発したシステムの成果を評価する。

## 2. 関連技術

### (1) 自然言語処理

自然言語処理とは、コンピュータが人間の言語を理解、解釈、捜査できるようにする技術である。人間が書いた文章から、感情や意見を抽出するなどの研究が行われている<sup>1)</sup>。自然言語処理では、人間のコミュニケーションとコンピュータ側の理解との間にあるギャップを解消することを目指している。本研究では、レビューの内容を自然言語処理にて解析して、参考になる度合いを測り、参考になる順番に並び替えを行う。

### (2) 正規化

正規化とは、データなどがある基準や形式に適合するように、一定の手順や規則に従って変形・変換することである。様々な分野で用いられる概念であり、それぞれ目的や方法などは異なる。自然言語処理において、どのように正規化するのがよいか、その方法についても研究が行われている<sup>2)</sup>。自然言語処理における正規化として、文字の種類統一、つづりや表記ゆれの吸収といった単語を置き換えるという方法がある。自然言語は書いた人によって様々な表記があり、同じ意味であるのに異なった言葉が使われることが多い。正規化して同じ意味を表す言葉を同じ表記とすることで、機械学習の際に計算量やメモリ使用量の軽減をすることが可能となる。例えば、自然言語処理において数字自体にはあまり意味を持たないことがほとんどであるため、数字は統一した単語(0など)に置き換えるという方法がある。また、英語のアルファベットの太文字と小文字の違いや、日本語のひらがな表記、カタカナ表記、漢字の違いは、同じ意味を表すので、統一するという方法がある。

### (3) ストップワード

ストップワードとは、検索の対象外とする言葉のリストのことである。自然言語処理において、全ての言葉を解析しようとする膨大な情報量となり、コンピュータではメモリ容量不足となったり、計算に長い時

間がかかったりと問題が発生する。そこで、情報検索に関係がない言葉や情報量が少ない言葉をストップワードとして定義し、検索の対象外とすることで、コンピュータで処理できる程度の情報量で自然言語処理を行う。例えば、日本語における「は」や「です」などの助詞や助動詞などの付属語、英語における be, have のような一般的な意味を持つ語など、情報検索に有効でないと考えられる語はストップワードとすることが多い。ストップワードをどのように定義するか、ということも自然言語処理では重要な検討事項である<sup>3)</sup>。

### (4) 学習素性

学習素性(feature for machine learning)とは、機械学習の際、データを分類する手がかりとなる情報のことである。単に「素性」ともいう。例えば、“単語の品詞を推定する”モデルを機械学習する際には、“その単語の前後に出現する語や品詞”に着目することで学習できると考えられるため、その情報が学習素性として使われる<sup>4)</sup>。つまり、前後に出現する語や品詞を手がかりとすることで対象語の品詞を推定するモデルを学習することができる。学習素性として何を使うかは、機械学習に基づく自然言語処理の成否を決める重要な要因である。本研究においては、レビューの文章に登場した頻度に応じて学習素性を決定する。

### (5) 機械学習

機械学習(machine learning)とは、データの集合から、未知のデータをあらかじめ定義されたいくつかのカテゴリに分類するモデルを自動的に学習する手法である。正解のカテゴリが付与されたデータを訓練データとする教師あり学習(supervised learning)と、正解のカテゴリが付与されていないデータを訓練データとする教師なし学習(unsupervised learning)が存在する。自然言語処理においても機械学習を用いる手法が用いられている<sup>5)</sup>。機械学習のひとつとして、本研究ではロジスティック回帰<sup>6)</sup>を用いる。ロジスティック回帰とは、多変量解析の手法のひとつで、目的変数が2つ以上ある場合に用いられる。自然言語処理においてもロジスティック回帰を利用した研究が行われている<sup>7)</sup>。

### (6) MeCab

MeCabとは京都大学情報学研究科-日本電信電話株式会社コミュニケーション科学基礎研究所共同研究ユニットプロジェクトを通じて開発されたオープンソース形態素解析エンジンである<sup>8)</sup>。特徴として、言語、辞書、コーパス(テキストや発話を大規模に集めてデータベース化した言語資料)に依存しない汎用的な設計であること、他のライブラリに比べ高速であること

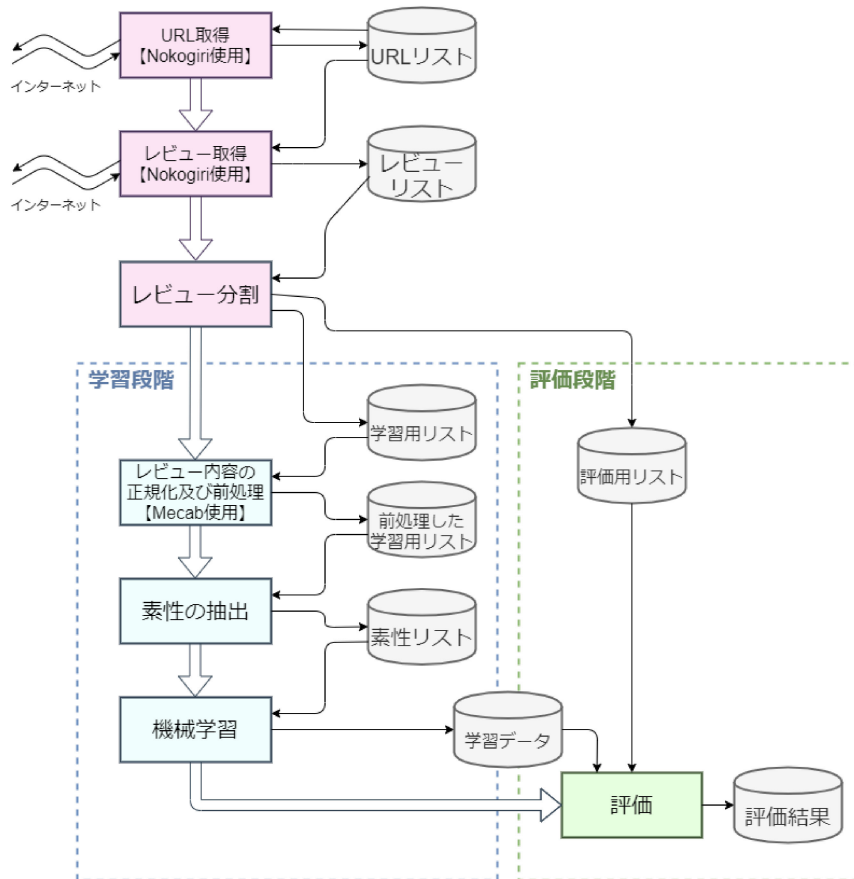


図1 システム概要図

などがあげられる。本研究では、レビューの文章の正規化など学習の前処理にて使用する。

### (7) Nokogiri

Nokogiri とは、Web スクレイピングを行うときによく使用される Ruby のライブラリである<sup>9)</sup>。Nokogiri を使用することで、HTML の解析や Web データの収集、すなわち Web スクレイピングを行うことができる。本研究では商品リンクの取得、レビューの取得のときに使用する。

## 3. 提案システム

### 3.1 概要

本研究で提案する、レビュー内容を評価して参考になる順番に並び替えを行うシステムの概要図を図1に示す。

提案システムに使用するレビューのデータはインターネット経由で Amazon から取得する。取得したレビューは、学習に用いるものと評価に用いるものに分割する。提案システムの処理は学習段階と評価段階のふたつに分けられる。学習段階では、インターネットか

ら取得したレビューに対して、正規化などの前処理を行い、機械学習で学習を行う。評価段階では、学習に使用していないレビューに対して、そのレビューが参考になると判断した順番に並べ替えを行い、その並べ方について評価値を算出する。

### 3.2 レビューの取得と分割

提案システムで使用するレビューデータを取得するために、まず商品リンクの URL の取得を行う。商品リンクの URL 取得のフローを図2に示す。

Amazon のページにはいくつかの商品へのリンクが貼られているので、それら商品リンクを Ruby のライブラリである Nokogiri を使用して Web スクレイピングを行う。Amazon のページ内の商品リンクは以下のように HTML で記述されている。

```
<a class="a-link-normal" href="商品リンクのURL">
```

Nokogiri を用いて、タグ a で属性 class が "a-link-normal" であるものを抽出すれば、属性 href に記述されている商品リンクの URL を取得することができる。更に URL を辿って商品ページの解析を繰り返すことで、多くの商品ページから URL を取得することができる。

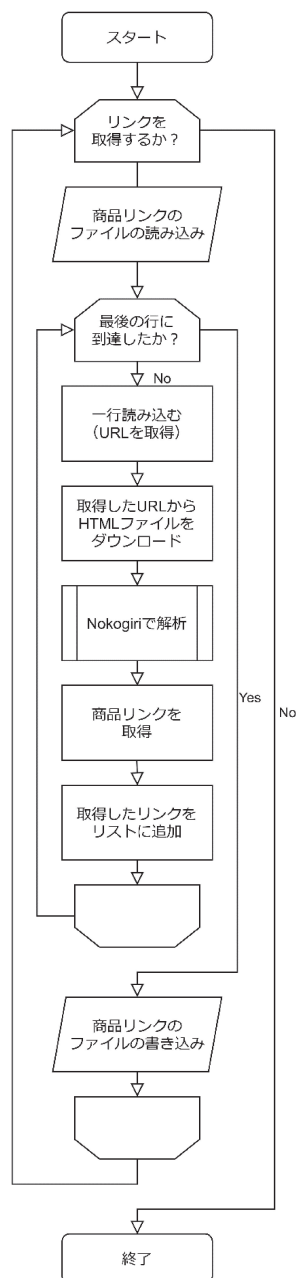


図2 商品リンクの URL 取得フロー

取得した URL は URL リストとしてテキストファイルで保存しておく。

次に、URL リストに保存している商品リンクのページにアクセスして、レビューの取得を行う。レビュー取得フローを図3に示す。

Amazon の商品ページには、商品名や値段の他に、レビューの数や評価値（5段階評価）が含まれている。商品リンクの URL と同様に、HTML に含まれるタグや属性を検出して、そのレビューの情報を取得する。例えばカメラという商品名の場合のは以下のように HTML で記述されている。

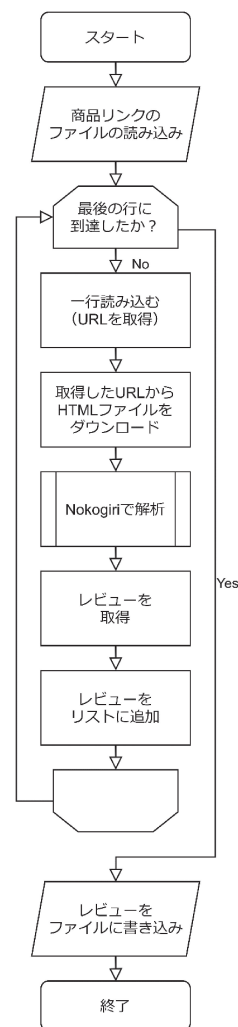


図3 レビュー取得フロー

<span id="productTitle" class="a-size-large product-title-word-break">

カメラ

</span>

従って、タグ span の属性 id="productTitle" を検出することで、商品名を取得することができる。取得したレビューはレビューリストとしてテキストファイルで保存しておく。取得するレビューの情報と、対応する HTML タグ・属性を表1に示す。

さらに、レビューリストに保存したレビューは、学習用リストと評価用リストの2つに分割する。学習用リストは学習段階で行う機械学習の教師データとして用いる。評価用リストは、機械学習にとって未知のデータであるため、評価に用いることができる。本研究では Web スクレイピングにより取得できた商品リンクの中で、5つのカテゴリの合計56件の商品リンクを使用した。その56件の商品に対して投稿された合



表1 取得するレビューの情報と対応する HTML タグ・属性

情報	HTML タグ	属性
商品名	span	id="productTitle"
商品の種類	a	class="a-link-normal a-color-tertiary"
商品の値段	span	id="priceblock_ourprice"
商品の評価数	span	id="acrCustomerReviewText"
質問の回答数	a	id="askATFLink"
レビュー	div	id="cm_cr-review_list"
5段階評価値	i	data-hook="review-star-rating"
レビュータイトル	a	data-hook="review-title"
プロフィール名	span	class="a-profile-name"
コメント日	span	data-hook="review-date"
コメント内容	span	data-hook="review-body"
参考になったと 考える人数	span	data-hook="helpful-vote-statement"

表2 商品のカテゴリ、商品数およびレビュー数

カテゴリ	商品数 [件]	レビュー数 [件]
Kindle ストア	22	7663
パソコン・周辺機器	5	11317
家電&カメラ	8	24921
本	19	3329
Amazon デバイス	2	5173
合計	56	52403

計 52,403 件のレビューを使用した。使用した商品のカテゴリ、商品数およびレビュー数について、表 2 に示す。

このレビューからランダムに選択した 1000 件を評価用リストとして保存し、残り 51,403 件を機械学習の教師データとするために学習用リストとして保存した。

### 3.3 学習段階

#### 3.3.1 教師データの準備処理

学習段階においては、最初に学習用リストに含まれるレビュー内容の正規化と前処理、素性の抽出を行う。これにより学習するための教師データを準備する。レビュー内容の正規化と前処理の手順を図 4 (a) に示す。

各正規化を行った後、MeCab で形態素の解析を行う。解析によって動詞や助詞など品詞ごとに分割することができる。動詞であれば活用語の変換行い、また、助詞などのストップワードは除去する。このような前処理を行うことで、表記ゆれをなくし、学習と関連が強い部分だけ残すことができる。各処理でどのような変換がおこなわれるか具体的な例文で以下に示す。変換が行われた部分は下線で示す。変換の内容は【 】内にて説明する。

##### (1) 正規化前

ロボットは 2 4 時間働けるので A I に仕事をとられる。

##### (2) Unicode 正規化

ロボットは 24 時間働けるので AI に仕事をとられる。【全角を半角に変換】

##### (3) 数字の正規化

ロボットは 0 時間働けるので AI に仕事をとられる。【数字を 0 に変換】

##### (4) 英字の正規化

ロボットは 0 時間働けるので ai に仕事をとられる。【英字は小文字に変換】

##### (5) Mecab で解析し、日本語の形態素を解析

ロボット\_は\_0\_時間\_働ける\_ので\_ai\_に\_仕事\_をとら\_れる。【品詞ごとに“\_”で分割】

##### (6) 活用語の原型への変換

ロボット\_は\_0\_時間\_働ける\_ので\_ai\_に\_仕事\_を\_とる\_れる。【動詞を原型に変換】

##### (7) ストップワードの除外【括弧で示す語は除去する】

ロボット\_(は)\_0\_(時間)\_働ける\_(の)\_ai\_(に)\_仕事\_(を)\_とる\_(れる)。

##### (8) 正規化後

ロボット\_0\_働ける\_ai\_仕事\_とる。

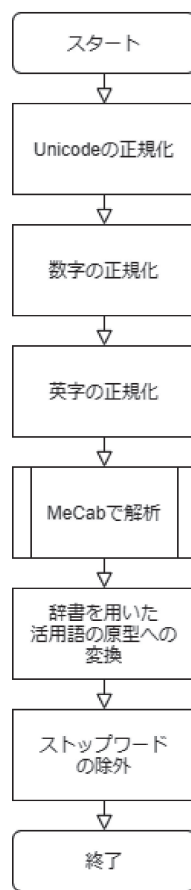
このように正規化やストップワードの除外などを行うことによって、学習に有意義な情報だけを残すことができ、コンピュータで計算ができる情報量とすることが可能となる。また、A I (全角大文字) や AI (半角大文字) など表記ゆれは ai (半角小文字) に統一されることになるので、学習するときに同じ語と認識することが可能となる。

正規化と前処理を行った後、素性の抽出を行う。本研究はレビューが参考になる順序に並べ替えることが目的であるため、レビューの文章に登場する回数に着目する。例えば「おいしい」という言葉がレビューに多く記入されていれば、「おいしい」という言葉はレビューにおいて重要な情報であると考えられる。そこで本研究では、レビューの文章に登場する回数に応じて、素性を決定する。素性の抽出のフローチャートを図 4 (b) に示す。

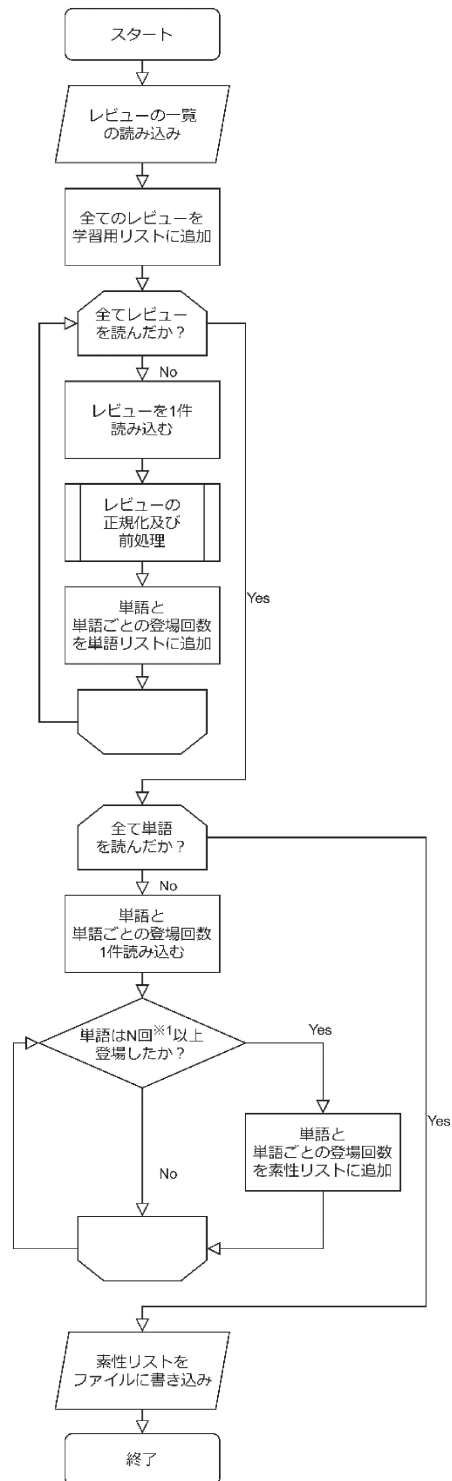
#### 3.3.2 機械学習

教師データの準備ができたなら、次に機械学習を用いてレビューの学習を行う。使用するモデルについて、線形回帰やロジスティック回帰が考えられるが、本研究は 0 から 1 までの実数となる 1 次元の評価値で表すことから、ロジスティック回帰を採用した。

目的変数をレビューが参考になる確率  $P$  とし、説明変数を素性  $x_i$  としたとき、ロジスティック回帰モデルは式 (1) で示される。



(a) レビューの正規化及び前処理



※1 Nは素性として抽出されるための  
単語の最低登場回数

(b) 素性の抽出

図4 教師データの準備処理のフロー

$$P = \frac{1}{1 + \exp(\theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_N x_N)} \quad (1)$$

ここで  $\theta_i$  はそれぞれの素性の重みを表す。N は素

性の数を表す。つまり、レビューが参考になる確率  $P$  は、レビューの内容に登場する回数で決定される N 個の素性  $x_i$  が、どの程度レビュー内容に影響するか

重み  $\theta_i$  で決定される。ただし、重み  $\theta_i$  は未知の値であるため、機械学習で適切な  $\theta_i$  を求める。なお、素性を含まないレビューも存在するため、 $x_0$  は“学習素性なし”として常に 1 とする。式 (1) はシグモイド関数とよばれ、確率  $P$  は  $(-\infty, \infty) \rightarrow (0, 1)$  の単調増加連続関数となり、1 つの変曲点を持つ。

ロジスティック回帰では、予測値  $P^*$  と正解の値  $P$  の差から損失関数  $L$  を求め、 $L$  が最小となる回帰直線を求める。回帰直線を求めるために最急降下法を用いた。最急降下法による損失関数  $L$  の最小値を求めるための勾配は式 (2) にて算出する。

$$\frac{\partial L}{\partial \theta} = - \sum_{n=1}^N x_n (P - P^*) \quad (2)$$

また、重み  $\theta_i$  の更新は式 (3) にて計算する。ここで  $\eta$  は重み  $\theta_i$  を更新する度合いを表すパラメータで、学習率という。

$$\theta \leftarrow \theta - \eta \frac{\partial L}{\partial \theta} \quad (3)$$

これにより、予測値  $P^*$  と正解の値  $P$  の差が最も小さくなるように重み  $\theta_i$  を更新することで学習を行う。学習した結果は、学習データとして保存する。

### 3.4 評価段階

機械学習によって学習データが作成できたら、分割しておいた評価用リストをその学習データを用いて評価を行う。評価については 4 章で説明する。

## 4. 評価

### 4.1 概要

本研究の目的はレビューを参考になる順序に並び替えることである。提案システムの評価は、実際の商品ページの並びと一致率で評価する。実際の商品ページに備えてある評価値による並べ方を  $L_R$  とし、これを正解の並べ方とする。評価値を使用せず、本研究の提案システムによる並べ方を  $L_P$  とし、 $L_R$  と  $L_P$  のそれぞれの要素の一致率を  $P_{\text{match}}$  とする。例えば、レビュー数 5 件の  $L_R$  と  $L_P$  が以下のように一致した場合、 $P_{\text{match}} = 100\%$  となり最も良い結果となる。

$L_R : \{ 1, 2, 3, 4, 5 \}$

$L_P : \{ 1, 2, 3, 4, 5 \}$

また、以下の場合にはレビュー 5 件中 4 件が一致しているので、 $P_{\text{match}} = 3/5 = 60\%$  となる。

$L_R : \{ 1, 2, 3, 4, 5 \}$

$L_P : \{ 1, 2, \underline{4}, \underline{3}, 5 \}$

ただし、このような単純な一致率による評価には欠

点がある。以下の場合にはレビュー 5 件中 1 件しか一致していないため、 $P_{\text{match}} = 1/5 = 20\%$  となる

$L_R : \{ 1, 2, 3, 4, 5 \}$

$L_P : \{ \underline{4}, \underline{1}, \underline{2}, \underline{3}, 5 \}$

しかし、この場合は 4 番目のレビューだけが一致しなかっただけであり、他のレビューの並べ方は正しいものが求められていると考えられる。従って、並べ方の評価としては良い結果であると考えられるが、一致率による評価の場合は 20% という低い評価値になってしまう。このような欠点を解消すべく、本研究ではクイックソートを利用した評価法 (Quicksort Evaluation method: QE 法) を提案する。

### 4.2 QE 法

並べ替えの評価として、クイックソート<sup>10)</sup> を利用した評価法である QE 法を提案する。クイックソートとは、データの並べ替え方法のひとつで、最も高速な方法といわれている。数値を並べたリストにおいて、クイックソートは以下のようなアルゴリズムである。

1. リストから基準となる値 (ピボット) を選択する
2. ピボットより小さい数をリストの前方へ、ピボットより大きい数をリストの後方に移動させる
3. ピボットの前方と後方に、それぞれ再帰的にクイックソートを適用する

本研究で提案する QE 法では、ピボットをリストの中央値とする。QE 法による評価値  $P_{QE}$  は、リストの要素を入れ替えた数  $S_{\text{count}}$  と、要素数における最も多い入れ替え回数  $S_{\text{max}}$  を用いて式 (4) で示される。

$$P_{QE} = 1 - \frac{S_{\text{count}}}{S_{\text{max}}} \quad (4)$$

$S_{\text{max}}$  はその要素数における最も多い入れ替え回数であるため、全ての要素が逆順になっている状態から、昇順になるようにクイックソートを行うときの入れ替え回数である。要素数  $n$  における  $S_{\text{max}}$  を  $S_{\text{max}}(n)$  で表す。例えば  $S_{\text{max}}(3)$  は以下のように考えられる

Step1. 全ての要素数が逆順になっている状態。ピボットは 2 とする

$\{ 3, 2, 1 \}$

Step2. クイックソートによる入れ替え。ピボットより小さい値 (1) をリスト前方に移動

$\{ 1, 3, 2 \}$

Step3. クイックソートによる入れ替え。ピボットより大きい値 (3) をリスト後方に移動

$\{ 1, 2, 3 \}$

以上で昇順に並び替えられるため、 $S_{\text{max}}(3) = 2$  とな

る. なお,  $S_{\max}(n)$  は式 (5) で表すことができる.

$$S_{\max}(n) = \begin{cases} 0 & (n = 0, 1) \\ n - 1 + S_{\max}\left(\frac{n}{2}\right) + S_{\max}\left(\frac{n}{2} - 1\right) & (n \geq 2 \text{ かつ偶数}) \\ n - 1 + 2 \cdot S_{\max}\left(\frac{n-1}{2}\right) & (n \geq 3 \text{ かつ奇数}) \end{cases} \quad (5)$$

リストが単純に一致するかを示す一致率  $P_{\text{match}}$  と, QE 法による評価値  $P_{\text{QE}}$  を比較する. いま, 要素の数が 5 のリスト  $L = \{5, 1, 2, 3, 4\}$  について,  $L$  は 5 という数値以外は昇順に並べられているため, 並べ替えにおける評価値は直感的に高い値になることが予想される. しかし, 一致率だけで考えると,  $P_{\text{match}} = 1/5 = 0.2$  となり評価値が低くなってしまふ. 一方, QE 法を用いると評価値は式 (6) で算出できる.

$$P_{\text{QE}} = 1 - \frac{S_{\text{count}}}{S_{\max}(5)} = 1 - \frac{2}{6} = 0.67 \quad (6)$$

なお, ここでの  $S_{\text{count}}$  の求め方は以下手順となる.

Step1. 初期状態. ピボットは 3

$\{5, 1, 2, 3, 4\}$

Step2. ピボット (3) を含んだ小さい数 (1, 2, 3) を前方に移動

$\{1, 2, 3, 5, 4\}$

Step3. ピボット (3) より大きい数 (4, 5) を前方に移動

$\{1, 2, 3, 4, 5\}$

このように,  $P_{\text{match}}$  では 0.2 と低くなってしまふ場合でも, QE 法を用いると約 0.67 という比較的高い評価値を算出することができる. この QE 法で得られる評価値が, 本研究における最終的な評価結果となる. リストの並び方と  $P_{\text{match}}$  及び  $P_{\text{QE}}$  との対応例を表 3 に示す.

## 5. 実験結果

学習段階において, 51,403 件のレビューを学習させ

表 3 リストの並び方と  $P_{\text{match}}$  及び  $P_{\text{QE}}$  との対応例

並び方	説明	$P_{\text{match}}$	$P_{\text{QE}}$
$\{1, 2, 3, 4, 5, 6\}$	昇順 (正しい並び)	1.000	1.000
$\{1, \underline{3}, \underline{2}, 4, 5, 6\}$	2 つだけ降順	0.667	0.875
$\{1, 2, \underline{5}, \underline{4}, \underline{3}, 6\}$	3 つが降順	0.500	0.750
$\{\underline{6}, 1, 2, 3, 4, 5\}$	最後尾が先頭になっている	0.000	0.750
$\{\underline{6}, 2, 3, 4, 5, \underline{1}\}$	2 つが入れ替わっている	0.667	0.500
$\{\underline{6}, \underline{1}, \underline{4}, \underline{5}, \underline{3}, \underline{2}\}$	ばらばら	0.000	0.375
$\{\underline{6}, \underline{5}, \underline{4}, \underline{3}, \underline{2}, \underline{1}\}$	降順 (正しい並びと逆)	0.000	0.000

※太字下線で示す数字は, 正しい並びからずれている

たシステムの評価を行うために実験を行う. 学習の素性とするか決定するための登場回数  $F$  は 100 から 10,000 の間の 7 つの値とした, 例えば,  $F = 1000$  であれば, 学習するレビューの中に「おいしい」という言葉が 1000 回出現した場合は, 「おいしい」という言葉を素性とし, 説明変数  $x$  とすることを意味する. 学習率  $\eta$  は 1.0 から 2.0 とした. 学習回数は 10,000 回とした. 実験対象である評価用リストに含まれるレビューは 1000 件とした. 登場回数  $F$  と学習率  $\eta$  の組み合わせごとの評価値  $P_{\text{QE}}$  を表 4 に示す.

ここで, 例えば  $F = 100$  の場合の素性数は 2083 となっているが, これは学習用リストのレビュー内容において, 100 回以上登場した言葉は 2083 語あったことを示す. 表 4 において, 評価値  $P_{\text{QE}}$  が高いところは赤く, 中間は白く, 低いところは青く色付けしている.

結果として,  $F = 5,000$ , 学習率 1.7 の場合に,  $P_{\text{QE}}$  が最も高い値 0.814 となった. ただし, 登場回数  $F$  や学習率  $\eta$  を変えても実験を行ったパラメータの組み合わせでは,  $P_{\text{QE}}$  は 0.712 を下回ることはなかった. なお,  $P_{\text{QE}}$  の値が 0.7 程度となるのは, リストが  $\{1, 2, 3, \underline{6}, 4, 5\}$  のような場合で, このときは  $P_{\text{QE}} = 0.75$  となる. このように  $P_{\text{QE}}$  の値が 0.7 程度となるのは, リスト内の数値が 1 つ誤っている程度である. 表 3 にも示したように, 6 つのリストの場合で最後尾のデータが先頭になっている場合など, 並べ方が大きくずれている場合で,  $P_{\text{QE}}$  は 0.50 となる. 従って  $P_{\text{QE}}$  が最大で 0.814 となった本システムは, レビューを参考になる順序に並び替える方法として有効であることが示された.

## 6. まとめ

インターネット通販サイトの商品レビューを, 参考になる順序に並び替えるシステムを構築するために, ロジスティック回帰を使用した機械学習にて商品レビュー内容を学習する方法を構築した. 学習対象として, Amazon のホームページから 52,403 件の商品レビューを取得し, うち 51,403 件を学習させ, 残り 1000 件で実験を行った. 評価については, クイックソートを利用し, 要素を入れ替えた回数に応じて評価値を算出する QE 法にて評価した. 実験の結果, 提案システムでは  $P_{\text{QE}}$  が 0.814 となり, 高い精度で商品レビューの並び替えが行えることが確認できた.

本研究ではインターネット通販サイトの商品レビューを対象にした. 商品レビューという人間が自然言語で記入した文章から, 参考になるかならないかを判断



表4 登場回数と学習率の組み合わせごとの評価値  $P_{QE}$ 

登場回数 F	素性数 N	学習率 $\eta$										
		1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7	1.8	1.9	2.0
100	2083	0.788	0.783	0.800	0.806	0.796	0.766	0.750	0.762	0.769	0.747	0.768
200	2083	0.788	0.783	0.800	0.806	0.796	0.766	0.750	0.762	0.769	0.747	0.768
500	1472	0.774	0.748	0.781	0.782	0.769	0.769	0.763	0.751	0.787	0.774	0.726
1000	1058	0.728	0.794	0.746	0.781	0.758	0.813	0.792	0.762	0.795	0.776	0.784
2000	701	0.782	0.756	0.781	0.712	0.737	0.734	0.722	0.800	0.795	0.769	0.718
5000	363	0.759	0.773	0.774	0.801	0.764	0.763	0.772	0.814	0.757	0.765	0.755
10000	207	0.795	0.804	0.795	0.809	0.789	0.782	0.794	0.781	0.741	0.787	0.759

することができた。今後は商品レビュー以外の文章を対象として、その文章が参考になるかならないか、さらにはその文章が真実であるか虚偽であるかという判定を行うことが検討課題としてあげられる。その場合は、学習素性をどのように定義するか、ロジスティック回帰以外にも、どのようなモデルで学習を行えばよいかの検討が考えられる。

#### 参 考 文 献

- 1) 小林のぞみ, 乾健太郎, 松本裕治, 立石健二, 福島俊一: 意見抽出のための評価表現の収集, *自然言語処理*, **12-2**, 203/222 (2005)
- 2) 泉朋子, 今村賢治, 菊井玄一郎, 藤田篤, 佐藤理史: 正規化を指向した機能動詞表現の述部言い換え, *言語処理学会第15回年次大会発表論文集*, 264/267 (2009)
- 3) 國府久嗣, 山崎治子, 野坂政司: 内容推測に適したキーワード抽出のための日本語ストップワード, *日本感性工学*, **12-4**, 511/518 (2013)
- 4) 中川哲治, 工藤拓, 松本裕治: Support Vector Machine を用いた未知語の品詞推定, *情報処理学会研究報告自然言語処理*, **2001-9**, 77/82 (2001)
- 5) 春野雅彦: 機械学習の手法による自然言語処理, *音声言語情報処理*, **29-23**, 133/138 (1999)
- 6) 高田直樹: ロジスティック回帰分析結果の解釈・利用のための新手法, *PROVISION*, **53**, 71/77 (2007)
- 7) 亀甲博貴, 三輪誠, 鶴岡慶雅, 森信介, 近山隆: ロジスティック回帰による言語モデルを用いた将棋解説文の自動生成, *言語処理学会第20回年次大会発表論文集*, 943/946 (2014)
- 8) McCab: Yet Another Part-of-Speech and Morphological Ana-

lyzer ホームページ: <https://taku910.github.io/mecab/> (2019年10月2日閲覧)

- 9) Nokogiri ホームページ: <https://nokogiri.org/> (2019年10月2日閲覧)
- 10) C. A. R. Hoare: Quicksort, *Comput. J.*, **5-4**, 10/15 (1986)

#### 著 者 紹 介

##### 市川 知春

2020年鹿児島工業高等専門学校卒業。2020年より鹿児島工業高等専門学校専攻科電気情報システム工学専攻に在学中。自然言語処理に関する研究に従事。

##### 武田 和大

2000年鹿児島大学大学院理工学研究科情報工学専攻（博士前期）修了。2006年鹿児島大学大学院理工学研究科情報工学専攻（博士後期）修了。博士（工学）。2000年鹿児島大学文部技官。2009年鹿児島工業高等専門学校情報工学科助教。2012年准教授。現在に至る。分散並列処理、群知能、気象データ、建築環境工学、電子物性に関する研究に従事。電気学会、空気調和・衛生工学学会、日本建築学会、情報文化学会会員。

##### 原 崇（正会員）

2003年鹿児島大学大学院理工学研究科情報工学専攻（博士前期）修了。2006年鹿児島大学大学院理工学研究科情報工学専攻（博士後期）修了。博士（工学）。2004年NECモバイルリング株式会社。2006年日本電気通信システム株式会社に勤務。2013年鹿児島工業高等専門学校情報工学科助教。2021年准教授。現在に至る。分散並列処理、群知能、自然言語処理に関する研究に従事。情報文化学会、空気調和・衛生工学学会、日本シミュレーション学会会員。