

フォーラムを教師データとしたアプリケーションレビュー分類手法の提案

市川 直人[†] 裕本 真佑[†] 楠本 真二[†]

[†] 大阪大学大学院情報科学研究科

E-mail: †{n-itikaw,shinsuke,kusumoto}@ist.osaka-u.ac.jp

あらまし ソフトウェア開発において、アプリケーションレビューは開発者にとって豊富な情報源である。しかし、レビューは膨大な数が投稿されており、無意味なレビューも多い。これらの情報を開発者が効率よく取得するために、自然言語処理と機械学習を用いてレビューを自動分類する手法が数多く報告されている。これらの手法では主に教師あり学習を用いるため、大量のレビューを目視で確認し、バグ報告や機能要求といったラベルを付与することで教師データを生成する必要がある。またレビューとは別に、エンドユーザが情報を共有する場としてフォーラムがある。特にゲームの分野で盛んな文化で、フォーラムに投稿される各トピックはバグ報告や機能要求などのカテゴリーに分類されている。本研究では、すでにカテゴライズされたフォーラムのトピックを教師データとして用いることで、教師データを用意するコストを削減したレビューの分類手法を提案する。実験では、分類精度は既存手法には劣るものの、提案手法は十分な精度でレビューを分類できることを示した。

キーワード フォーラム, アプリケーションレビュー, 機械学習, 自然言語処理, BERT

1 はじめに

ソフトウェア開発において、アプリケーションレビューは開発者にとって豊富な情報源である [1]。有益な情報の例としては、バグの報告や新機能の提案、機能強化の要望、特定の機能に対するユーザのフィードバックなどが挙げられる。

これらのレビューは、アプリの開発やメンテナンスに活用できる。しかし、レビューの中には開発者ではなくユーザを対象とした内容も数多く存在しており、これらはソフトウェアの改善という目的にはあまり役に立たない。例えば、good や bad などの数単語のみで構成される内容や、インストールできないなど説明の不足した問題の報告、顔文字のみの内容などが存在する。これらは対象アプリの印象を共有する、というレビューそのものが持つ目的に従う内容ではあるが、開発の改善には繋がらない。人気のアプリに対しては1日に数百件のレビューが投稿されているため、開発者による手作業での確認には限界がある。

そこで、膨大な数のレビューから開発者が有用な情報を効率よく取得することを目的とした研究が多く行われている [2-8]。これらの研究では自然言語処理と機械学習、主に教師あり学習を用いて、レビューをその内容ごとに自動で分類する。既存研究では主に、Apple App Store や Google Play のようなモバイルアプリ配信プラットフォームに投稿されたレビューを、バグ報告や機能要求といったカテゴリーに分類している。

既存研究の課題として、教師データを用意するコストの高さが挙げられる。機械学習において教師データの数と質は精度に大きく寄与するため、既存手法では大量のレビューを目視で精読し、バグ報告や機能要求といったラベルを付ける必要がある。

実際に、Maalej ら [4] は 4,400 件、Zhang ら [5] は 3,092 件のレビューを目視で確認しており、大きな労力を費やしている。

そこで本研究では、アプリレビュー分類における教師データ作成の労力削減を目的として、レビューを目視でラベル付けする代わりに、すでにカテゴライズされたフォーラム上のトピックを教師データとして用いる手法を提案する。フォーラムとは、ユーザと開発者が自由に議論できるインターネットコミュニティの一形態であり、個々のソフトウェアに対して設置されることが多い。例えば Web 会議サービスである Zoom には、Zoom Developer Forum^(注1) が設けられている。フォーラムはレビューとは異なり、バグ報告や機能要求、一般議論などのカテゴリーがフォーラム管理者によって事前に設けられている。よってフォーラムへの投稿(トピック)は、すでにカテゴリー毎に分類された自然言語のデータであると見なすことができる。このトピック群をレビュー分類の教師データとすることで、目視によるラベリングコストの大幅な削減が可能となる。

具体的な手法としては、アプリレビューをバグ報告、機能要求、その他の3種類に分類するために、フォーラムのバグ報告、機能要求、一般議論のカテゴリーに属するトピックをそれぞれに対応する教師データとして利用する。フォーラムが設けられているアプリはあまり多くないが、ゲームの分野ではフォーラムが設けられることが多い。そこで実験では、デジタルゲーム配信プラットフォームである Steam 上のフォーラムを持つゲームへのレビューを対象として、提案手法と既存手法の分類精度を比較した。その結果、提案手法は既存手法に比べて分類精度は低下する

(注1) : <https://devforum.zoom.us>

が、十分な精度でレビューを分類可能だと示した。

2 準備

2.1 既存手法の流れ

近年、膨大な量のアプリレビューの中から開発者にとって有益な情報を効率的に取得するために、機械学習を用いてアプリレビューを自動分類する研究が多く行われている [2-8]。これらの研究は、図 1 に示したフローに従ってレビューの分類を行う。

まずはデータの用意を行う。教師あり学習に用いる教師データを作成するため、アプリストアからいくつかのレビューを抽出し、目視によって個々のレビューがどの種類に属するかをラベル付ける。既存研究では主に、Apple App Store や Google Play のようなモバイルアプリ配信プラットフォームに投稿されたレビューを対象に実験を行っている。分類の区分は研究によって様々である。Chen ら [2] は、レビューを開発者にとって有益かどうかで 2 種類に分類し、Maalej ら [4] は、バグ報告、機能要求、ユーザ体験、評価の 4 種類に分類している。さらに Zhang ら [5] は、バグ報告や機能要求を含む計 17 種類のカテゴリに分類している。

次に自然言語処理を行う。自然言語処理とは、曖昧な性質を持つ自然言語をコンピュータが扱えるように加工する処理である。具体的には、小文字への統一化や単語単位への分割などの前処理や、文章をその特徴を表す分散表現（ベクトル）に変換する処理を行う。

続いて分類モデルの構築を行う。分類モデルに対して自然言語処理を施した教師データを与えることで、正しく文章を分類できるように学習を行う。分類モデルには、ナイーブベイズやロジスティック回帰のような古典的な機械学習手法 [2] [3] [4] や、深層学習を利用した手法 [8] など様々な手法が報告されている。

最後にレビューの分類を行う。学習済みモデルに自然言語処理を施した未分類状態のレビューを入力することで、各レビューがどの種類に属するかを判定する。

2.2 課題

既存研究の課題として、データの用意にかかるコストの高さが挙げられる。既存研究はそれぞれ、自然言語処理や分類モデルに独自のアプローチを施すことで分類精度の向上を目指している。一方で、データの用意については共通して、数千件のレビューを目視でラベル付けすることで教師データを生成している。機械学習には多くの教師データが必要であり、目視で教師データを作成する労力は非常に大きい。特にバグ報告や機能要求に該当するレビューは、レビュー全体に占める割合が少なく、十分な数を確保するためには大量のレビューを目視する必要がある。実際に Maalej ら [4] は、ラベリングを行うための GUI ツールを作成することで、Apple App Store と Google Play から取得した計 4,400 件のレビューを目視でラベル付けしている。このような目視で教師データを作成する労力は、アプリの開発者が、自身が開発したアプリに対するレビューを分類しようとした際に大きな障壁になると考えられる。

2.3 フォーラム

フォーラムとは、レビューと同様にエンドユーザが感想や質

問を投稿し、開発者やユーザ同士で自由に議論や情報共有ができるインターネットコミュニティの一形態である。Arduino Forum^(注2)や Zoom Developer Forum のように、一般にフォーラムは 1 つのアプリケーションやプロジェクト単位で設けられる。様々な内容の投稿が混在しているレビューとは異なり、フォーラムにはバグ報告、機能要求、一般議論のようなカテゴリがフォーラムをの管理者によって事前に設けられており、フォーラムへの投稿（トピック）はその内容毎に分類されている。これによって、開発者が開発に有益なトピックを確認しやすい仕組みが完成している。フォーラムはアプリの開発元が公式に運営していることもあれば、有志のエンドユーザが非公式に運営しているものもある。また、1 つのアプリに複数のフォーラムが存在することもあり、その規模や設けられているカテゴリは様々である。ただし、バグ報告や機能要求、一般議論に該当するカテゴリは多くの場合で共通して設けられている。フォーラムは全てのアプリに設置されているわけではなく、フォーラムを持つアプリはあまり多くない。しかし、アプリの中でもゲームの分野ではフォーラムを持つアプリが比較的多い傾向にある。

2.4 Steam

Steam^(注3)は、Valve Corporation が開発したデジタルゲーム配信プラットフォームである。30,000 以上のゲームが配信され、1 億 2,000 万人以上のアクティブユーザーがいるなど、PC ゲームのデジタル配信プラットフォームとしては最大級の規模を誇る。

Steam では、ユーザはプレイしたゲームのレビューを投稿することができ、1,000 万件以上のレビューが投稿されている。特に人気の高いゲームには 1 日に 500 件以上のレビューが投稿され、開発者がそのすべてに目を通すことは困難であることが示されている [9]。これらのレビューには Apple App Store や Google Play に代表される一般的なモバイルアプリ配信プラットフォームに投稿されるレビューと同様に、雑多な感想やアスキーアートのような開発者にとってあまり役に立たないレビューと、バグ報告や機能要求のような開発者にとって有益なレビューが混在している。

3 提案手法

本研究の目的は、アプリレビュー分類における教師データ作成の労力削減にある。そこで本研究では、目視でのラベル付けによって教師データを用意する代わりに、すでにカテゴライズされたフォーラムのトピックを教師データとする。これにより、目視でのラベル付けによる教師データ作成のコストを大幅に削減したレビュー分類手法を提案する。具体的には、レビューをバグ報告、機能要求、その他の 3 種類に分類するために、フォーラムのバグ報告、機能要求、一般議論のカテゴリに属するトピックをそれぞれに対応する教師データとして利用する。

図 2 に提案手法の流れを示す。既存手法との大きな違いは、データの用意のステップである。既存手法ではレビューをいくつか無作為に抽出して目視でラベル付けを行っているのに対し

(注2) : <https://forum.arduino.cc>

(注3) : <https://store.steampowered.com/>

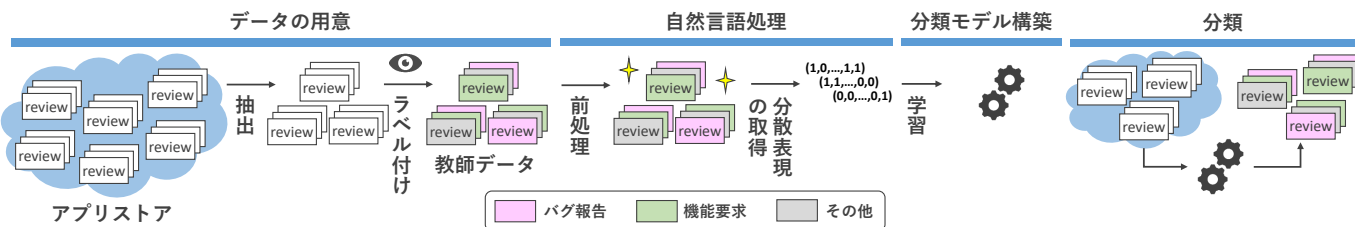


図1 既存手法の流れ

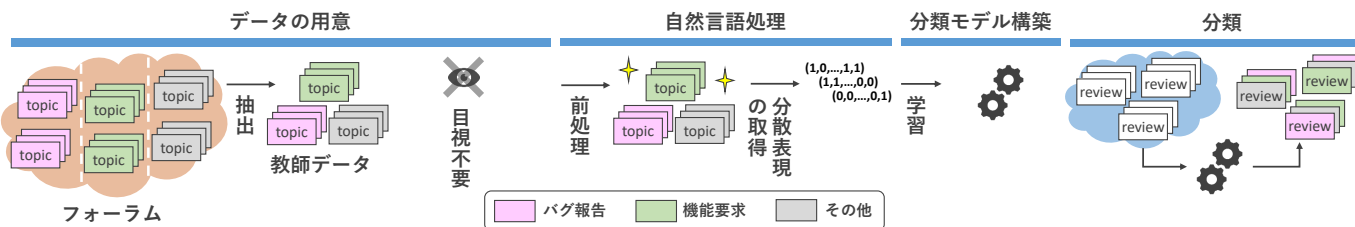


図2 提案手法の流れ

て、提案手法ではすでにカテゴリ化されたフォーラムのトピックを用いるため、目視によってラベル付けを行う必要がない。自然言語処理や分類モデルの構築、実際にレビューを分類するステップは既存手法と同様に行い、自然言語処理の技術や分類モデルは既存研究で用いられている技術の中でも特に分類精度の高かった技術を用いる。詳細は5.2節、5.3節にて述べる。

提案手法の最大の目的は、教師データを用意する労力の削減であり、分類精度の向上ではない。レビューとフォーラムはどちらもエンドユーザが自身の意見を投稿することのできる場であるが、一般にレビューはライトユーザからの投稿が多く、フォーラムはヘビーユーザからの投稿が多い傾向にある。そのため、レビューとフォーラムに投稿される文章には少なからず性質の違いが生じる。これにより、提案手法の分類精度は既存手法よりも低下することが予想される。一方で、提案手法では大量の教師データを用意することができるため、機械学習において有利に働く可能性もある。

4 Research Question

本研究では、提案手法の有効性を確認するために、以下の2つのResearch Question (RQ) を設ける。

RQ1: 提案手法の分類精度は既存手法と比較してどの程度か？

RQ1では、フォーラムのトピックが、レビューを分類する際の教師データとして機能するかどうかを確認する。ここで、あるアプリ app のフォーラムから得られたトピック群を F_{app} 、アプリストアから得られたレビュー群を R_{app} と定義する。さらに、 F_{app} を教師データ、 R_{app} をテストデータとする提案手法を、 $E(F_{app} \Rightarrow R_{app})$ と表記する。既存手法は、教師データとテストデータの両方にレビューを用いるため、 $E(R_{app} \Rightarrow R_{app})$ と表現できる。本RQでは、これら2つの手法の分類精度を比較する。3節で述べたように、レビューとフォーラムの性質の違いによって提案手法の分類精度は既存手法に比べて低下することが予想される。しかし、提案手法でも一定の精度が確認できれば、フォーラムのトピックが教師データとして機能することが確

認できる。

RQ2: フォーラムを持たないアプリへのレビューを分類する際にも提案手法を適用可能か？

2.3節で述べたように、フォーラムを持つアプリは多くはない。フォーラムを持たないアプリへのレビューを分類する際に提案手法を適用する場合、別アプリのフォーラムを教師データとして利用することが考えられる。異なるアプリのフォーラムを教師データとして用いる場合、そのアプリ固有の表現や文化が影響することで分類精度が低下する可能性がある。そこで、異なるアプリのフォーラムでも教師データとして機能するかを調査するために、あるアプリ app_A のフォーラムのトピックを教師データとして同じアプリへのレビューを分類する場合 $E(F_{app_A} \Rightarrow R_{app_A})$ と、別のアプリ app_B のフォーラムを教師データとしてあるアプリ app_A へのレビューを分類する場合 $E(F_{app_B} \Rightarrow R_{app_A})$ の分類精度を比較する。 $E(F_{app_B} \Rightarrow R_{app_A})$ でも一定の精度が確認できれば、フォーラムを持たないアプリへのレビューを分類する場合でも提案手法が適用できることが確認できる。

5 実験手順

5.1 データの用意

実験を行うにあたって、分類対象になるレビューと教師データになるフォーラムのトピックを収集する。2.3節で述べたように、ゲームにはフォーラムを設けられることが多い。そこで本実験では、ゲームに対するレビューを対象とし、Steam上からレビューを収集する。これは2.4節のとおり、Steamはゲーム配信プラットフォームとして最大級の規模を誇り、ゲームに対する多くのレビューが存在するためである。また、Steam上の全てのゲームがフォーラムを持っているわけではないため、本研究で用いるデータセットとして、Steamの売り上げ上位のゲームタイトルからフォーラムを持つタイトルを選定する。売り上げ上位のタイトルを対象とするのは、より多くのレビューを収集するためである。なお、本実験では開発元が公式に運営しているフォーラムを持つタイトルのみ対象とした。これらの条件を満た

表 1 実験データ

| | レビュー数 | | | | フォーラムのトピック数 | | | |
|------------------------|---------|--------|-------|-------|-------------|-------|-------|--------|
| | #バグ報告 | #機能要求 | #その他 | 合計 | #バグ報告 | #機能要求 | #一般議論 | 合計 |
| Cities: Skylines | 108(17) | 67(15) | 1,000 | 1,175 | 7,007 | 2,813 | 5,650 | 15,470 |
| Euro Truck Simulator 2 | 61(5) | 84(6) | 959 | 1,104 | 13,113 | 4,800 | 4,184 | 22,097 |

※括弧内の数値は、レビュー全体から無作為に抽出した約 1000 件に含まれるバグ報告、機能要求のラベルが付与されたレビューの数である。

すタイトルとして、Cities: Skylines^(注4)(Cities) と Euro Truck Simulator 2^(注5)(Euro) を本実験の対象として選定した。これらのタイトルのレビューのうち英語で書かれたものと、これらのタイトルのフォーラム^(注6)^(注7)からバグ報告、機能要求、一般議論に該当するカテゴリに分類されたトピックを収集した。

提案手法は、教師データを作成するために目視によるラベル付けを行わずにすむことがメリットであるが、本実験では、レビューを目視でラベル付けする作業を行う。これは、提案手法の分類精度の比較対象として既存手法を実施する必要があり、その既存手法の適用においてはレビューのラベル付けが必須であることが理由の一つである。さらに、提案手法と既存手法のいずれにおいても、その分類精度を確認するためには分類の正解集合を知る必要があることも理由として挙げられる。

収集したレビューからタイトルごとに約 1,000 件を無作為抽出し、目視によりバグ報告、機能要求、その他のラベルを付与した。しかし、バグ報告、機能要求に該当するレビューはそれぞれ 20 件以下であり、実験に十分な数が取得できなかった。そこでバグ報告によくみられる単語として“bug”, “fix”, “crash” [4]を含んだレビューの中から 800 件を無作為抽出し、目視によってラベル付けを行った。その中からバグ報告、機能要求のラベルが付与されたレビューを実験データに追加した。機能要求についても同様に、機能要求によくみられる単語として“please”, “hope”, “improve”, “need”, “prefer”, “request”, “suggest”, “wish” [4]を含んだレビューの中から 800 件を無作為抽出し、目視によってラベル付けを行った。その中からバグ報告、機能要求のラベルが付与されたレビューを実験データに追加した。

表 1 に収集したレビュー数とフォーラムのトピック数を示す。本研究ではレビューのラベル付け作業に 50 時間以上費やしたが、得られたレビュー数に対してフォーラムのトピック数は 10 倍以上ある。提案手法では教師データが大量に用意できることがわかる。

5.2 自然言語処理

自然言語処理技術を用いてテキストに前処理を施すことによって、自然言語の持つ曖昧さを削減し、分類精度が向上することが報告されている [4]。収集したレビュー、フォーラムに対して一般的な自然言語処理として、小文字化、Lemmatization、ストップワード除去を行った。

小文字化：大文字小文字の違いによる曖昧さを無くすために、

(注4) : https://store.steampowered.com/app/255710/Cities_Skylines/

(注5) : https://store.steampowered.com/app/227300/Euro_Truck_Simulator_2/

(注6) : <https://forum.paradoxplaza.com/forum/forums/cities-skylines.859/>

(注7) : <https://forum.scssoft.com/viewforum.php?f=3>

文章全体を小文字に統一する。

Lemmatization : Lemmatization とは、同一の語彙素をレンマ（その語彙素を代表する基本形）に統一する処理である。例えば、“fixing”, “fixed”, “fixes” は全て “fix” に統一される。

ストップワード除去：“the”, “am”, “their” のような、文法的な機能を持ち、文の意味にあまり影響を与えない一般的な単語をストップワードとして除去する。ストップワードは Maalej ら [4] によって定義されたものを用いた。

5.3 分類モデル構築

本実験で用いる分類器は BERT (Bidirectional Encoder Representations from Transformers) [10] を用いて構築する。BERT は 2018 年に Google から発表された自然言語処理を行うための深層学習モデルであり、転移学習によりさまざまなタスクに対応でき、少ないデータを追加で学習することによってモデルの構築を行うことができるという特徴がある。アプリレビューの分類においても、ナイーブベイズやロジスティック回帰による分類モデルに比べて BERT による分類モデルが高い精度を示したことが報告されている [11]。

本実験では、ラベル付けを行った Cities と Euro のレビューをそれぞれ 7:3 の比率で教師データとテストデータに分割し、 R_{Cities} , R_{Euro} , F_{Cities} , F_{Euro} の 4 種類の教師データでそれぞれ学習した 4 種類の分類器を生成した。なお、本実験では BooksCorpus と Wikipedia によって事前学習された BERT-Base モデルをファインチューニングして使用し、全ての学習はバッチサイズ 128, エポック数 5 で実施した。

5.4 評価

生成した分類器を用いて、 R_{Cities} と R_{Euro} を分類し、その精度を調査する。まず、既存手法に倣い $E(R_{\text{Cities}} \Rightarrow R_{\text{Cities}})$, $E(R_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ を実施する。次に、RQ1 を調査するために提案手法として $E(F_{\text{Cities}} \Rightarrow R_{\text{Cities}})$, $E(F_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ を実施する。最後に、RQ2 を調査するために、提案手法を交差的に適用した $E(F_{\text{Euro}} \Rightarrow R_{\text{Cities}})$, $E(F_{\text{Cities}} \Rightarrow R_{\text{Euro}})$ を実施する。

分類精度は、Precision (精度), Recall (再現率), F1-score, ROC 曲線の AUC (Area under the curve) の 4 つの評価指標を用いて評価する。ROC 曲線は、縦軸を真の陽性率、横軸を偽陽性率として、様々な閾値に対応する点をプロットした曲線で、AUC は ROC 曲線の下面積で定義される。AUC は分類器の性能を特定の閾値に拠らず評価する指標であり、0 から 1 の値を取る。AUC は分類器の性能が高いほど 1 に近づき、ランダムに分類を行う分類器に対しては 0.5 となる。一般に、AUC が 1.0~0.9 で高精度、0.9~0.7 で中精度、0.7~0.5 で低精度とされる [12]。

表 2 実験結果

| | バグ報告 | | | | 機能要求 | | | | その他 | | | |
|--|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | Prec. | Recall | F1 | AUC | Prec. | Recall | F1 | AUC | Prec. | Recall | F1 | AUC |
| $E(R_{\text{Cities}} \Rightarrow R_{\text{Cities}})$ | 0.86 | 0.86 | 0.86 | 0.98 | 0.38 | 0.30 | 0.33 | 0.93 | 0.95 | 0.96 | 0.96 | 0.97 |
| $E(F_{\text{Cities}} \Rightarrow R_{\text{Cities}})$ | 0.48 | 0.79 | 0.59 | 0.94 | 0.32 | 0.65 | 0.43 | 0.77 | 0.96 | 0.84 | 0.90 | 0.79 |
| $E(F_{\text{Euro}} \Rightarrow R_{\text{Cities}})$ | 0.31 | 0.18 | 0.23 | 0.77 | 0.29 | 0.40 | 0.33 | 0.85 | 0.89 | 0.90 | 0.90 | 0.77 |
| $E(R_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ | 0.75 | 0.43 | 0.55 | 0.98 | 0.35 | 0.79 | 0.48 | 0.94 | 0.97 | 0.89 | 0.93 | 0.95 |
| $E(F_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ | 0.30 | 0.57 | 0.39 | 0.77 | 0.30 | 0.67 | 0.42 | 0.90 | 0.93 | 0.78 | 0.85 | 0.81 |
| $E(F_{\text{Cities}} \Rightarrow R_{\text{Euro}})$ | 0.62 | 0.62 | 0.62 | 0.88 | 0.33 | 0.46 | 0.39 | 0.78 | 0.93 | 0.90 | 0.91 | 0.81 |

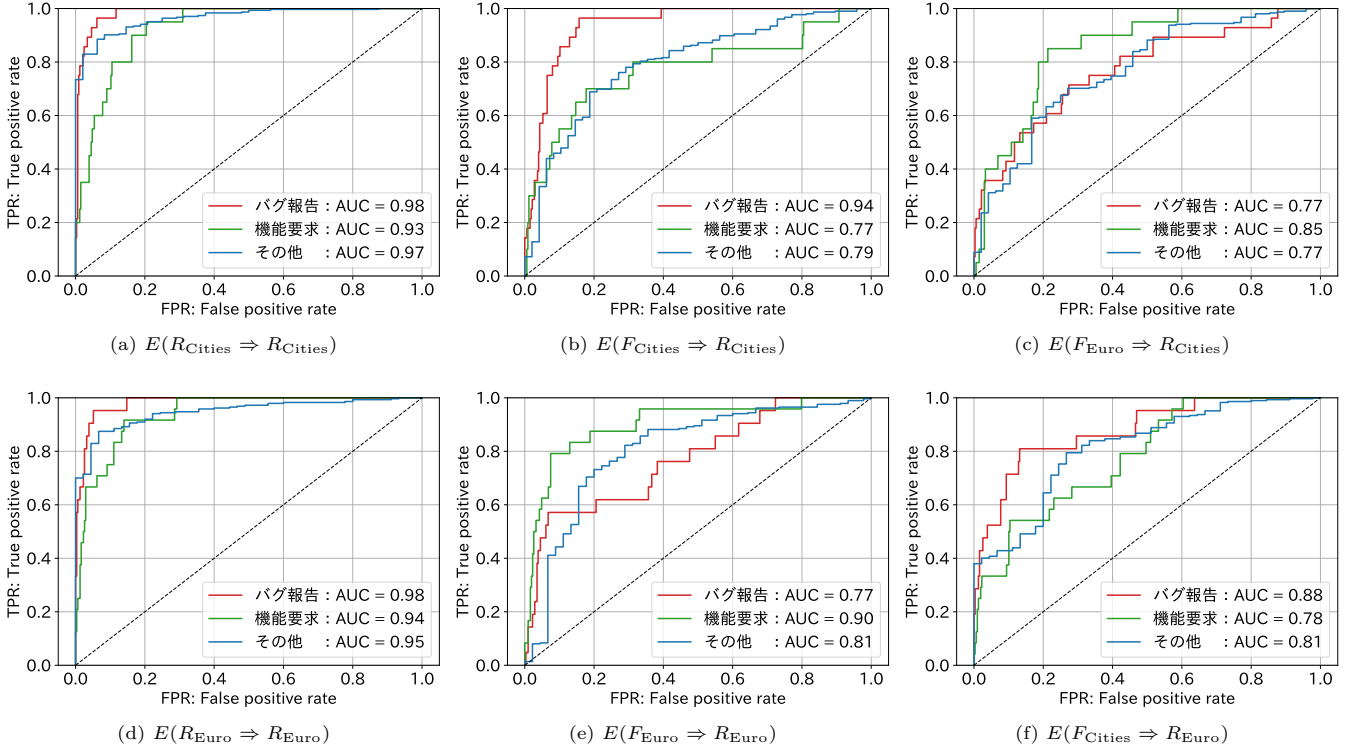


図 3 ROC 曲線

6 実験結果

本実験で実施した各手法の結果の Precision, Recall, F1-score, AUC を表 2 に示す。太字の数値は、それぞれのタイトルについて各列での最高値を意味している。また、表 2 の AUC を算出するための ROC 曲線を図 3 に示す。

6.1 RQ1 の結果

RQ1 では、既存手法と提案手法の比較として、 $E(R_{\text{Cities}} \Rightarrow R_{\text{Cities}})$ と $E(F_{\text{Cities}} \Rightarrow R_{\text{Cities}})$ の比較、 $E(R_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ と $E(F_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ の比較をそれぞれ行う。表 2 を見ると、既存手法は、どちらのタイトルでも全体的に高い精度を示している。特に AUC の値は全て 0.9 を超えており、非常に高い精度で分類ができていることがわかる。対して提案手法では、どちらのタイトルでも全体的に既存手法に劣る結果となったものの、AUC は全て 0.7 を超えており、十分な精度で分類ができていると言える。分類精度が低下した原因としては、4 節の通り、フォーラムとレビューの性質の違いにあると考えられる。また、5.3 節で述べたように、BERT は膨大なデータで事前学習したモデルを元

に、少ないデータを追加で学習してファインチューニングすることでモデルを構築できるという特徴があるため、提案手法の大量に教師データを用意できるという有利がうまく働かなかったことも原因の一つに挙げられる。

6.2 RQ2 の結果

RQ2 では、フォーラムを持たないアプリのレビューを分類するシナリオを想定して、提案手法を交差的に適用した場合の分類精度を確認する。図 3 の $E(F_{\text{Cities}} \Rightarrow R_{\text{Cities}})$ と $E(F_{\text{Euro}} \Rightarrow R_{\text{Cities}})$ を比較すると、バグ報告の AUC が減少して機能要求の AUC が増加している。対して、 $E(F_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ と $E(F_{\text{Cities}} \Rightarrow R_{\text{Euro}})$ を比較すると、バグ報告の AUC が増加して機能要求の AUC が減少している。つまり交差的に適用しているかどうかに関係なく、 F_{Cities} を教師データとして学習した分類器はバグ報告の AUC が高く、 F_{Euro} を教師データとして学習した分類器は機能要求の AUC が高いことが読み取れる。そこで、分類対象のレビューを固定して比較するのではなく、教師データを固定して比較を行う。 $E(F_{\text{Cities}} \Rightarrow R_{\text{Cities}})$ と $E(F_{\text{Cities}} \Rightarrow R_{\text{Euro}})$ を比較すると、バグ報告の AUC が多少減少しているものの、それ

以外に大きな差は見られない。同様に、 $E(F_{\text{Euro}} \Rightarrow R_{\text{Euro}})$ と $E(F_{\text{Euro}} \Rightarrow R_{\text{Cities}})$ を比較しても、機能要求の AUC が多少減少しているが、あまり大きな差は見られない。これらの結果から、アプリ固有の表現や文化は分類精度にあまり影響を与えないことがわかる。また、交差的に提案手法を適用した場合でも AUC は全て 0.7 を超えており、フォーラムを持たないアプリでも、他のアプリのフォーラムのトピックを教師データとして用いることで、十分な精度で分類ができると言える。

7 妥当性の脅威

データセットのラベル付けが分類器の性能に影響している可能性がある。本実験ではまず初めに、第一著者と第二著者の 2 人で数十件のレビューを確認し、バグ報告、機能要求の基準を定めたのちに、第一著者が 1 人でラベル付けを行った。そのため主観の影響が完全には排除できていない。

また、本実験で行った学習の各パラメータは筆者が経験則で定めたものであり、パラメータチューニングを行っていない。パラメータチューニングを行い、より適切なパラメータで実験を行うことで分類精度が向上する可能性がある。

本研究では、Cities: Skylines と Euro Truck Simulator 2 の 2 つのゲームのフォーラムを対象として実験を行った。しかし、フォーラムによってその規模やカテゴリーの区分は様々である。本実験で対象としたフォーラムとは別のフォーラムを教師データとした場合に、本実験と同程度の分類精度を達成できるかどうかは明らかでない。

8 おわりに

本研究では、フォーラムを用いることで、教師データを用意するコストを削減したアプリレビュー分類手法を提案した。実験では、デジタルゲーム配信プラットフォームである Steam 上のレビューを対象に、レビューをバグ報告、機能要求、その他の 3 種類に分類した。その結果、提案手法は既存手法と比べて精度は下がるものの、十分な精度で分類ができることを示した。またフォーラムがないアプリについても、別のアプリのフォーラムを教師データとすることで提案手法が適用できることを示した。

今後の課題としては、対象とするゲームタイトルを増やし、一般化可能性を向上させることが最優先の課題として挙げられる。また今回はゲームドメインを対象に実験を行ったが、既存研究で主に対象とされているモバイルアプリでの実験も考えられる。

謝辞 本研究は、日本学術振興会科学研究費補助金基盤研究 (B) (課題番号: 18H03222) の助成を得て行われた。

文 献

- [1] A. Holzer and J. Ondrus, "Mobile application market: A developer's perspective," *Telematics and Informatics*, vol.28, no.1, pp.22–31, 2011.
- [2] N. Chen, J. Lin, S.C.H. Hoi, X. Xiao, and B. Zhang, "AR-Miner: Mining Informative Reviews for Developers from Mobile App Marketplace," In Proc. International Conference on Software Engineering, pp.767–778, 2014.
- [3] E. Guzman, M. El-Haliby, and B. Bruegge, "Ensemble Methods for App Review Classification: An Approach for Software Evolution," In Proc. International Conference on Automated

Software Engineering, pp.771–776, 2015.

- [4] W. Maalej, Z. Kurtanović, H. Nabil, and C. Stanik, "On the Automatic Classification of App Reviews," *Requirements Engineering*, vol.21, no.3, pp.311–331, 2016.
- [5] L. Zhang, X.-Y. Huang, J. Jiang, and Y.-K. Hu, "CSLabel: An Approach for Labelling Mobile App Reviews," *Journal of Computer Science and Technology*, vol.32, pp.1076–1089, 2017.
- [6] C. Stanik, M. Häring, and W. Maalej, "Classifying Multilingual User Feedback using Traditional Machine Learning and Deep Learning," In Proc. International Requirements Engineering Conference Workshops, pp.220–226, 2019.
- [7] M.B. Messaoud, I. Jenhani, N.B. Jemaa, and M.W. Mkaouer, "A Multi-label Active Learning Approach for Mobile App User Review Classification," In Proc. Knowledge Science, Engineering and Management, pp.805–816, Springer International Publishing, Cham, 2019.
- [8] N. Aslam, W.Y. Ramay, K. Xia, and N. Sarwar, "Convolutional Neural Network-Based Classification of App Reviews," *IEEE Access*, vol.8, pp.185619–185628, 2020.
- [9] D. Lin, C.-P. Bezemer, Y. Zou, and A.E. Hassan, "An empirical study of game reviews on the Steam platform," *Empirical Software Engineering*, vol.24, no.1, pp.170–207, 2019.
- [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," In Proc. Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp.4171–4186, 2019.
- [11] 山田侑樹, 樋山淳雄, "BERT によるアプリケーションレビュー分類モデルの構築と評価," *電子情報通信学会技術研究報告*, vol.120, no.423, pp.25–30, 2021.
- [12] J. Fischer, L. Bachmann, and R. Jaeschke, "A readers' guide to the interpretation of diagnostic test properties: Clinical example of sepsis," *Intensive care medicine*, vol.29, pp.1043–1051, 2003.