

# OSS プロジェクトへのオンボーディング支援のための Good First Issue 自動分類

堀口 日向 大平 雅雄

オープンソースソフトウェア (OSS) 開発プロジェクトは、プロジェクトの持続可能性を維持するために常に新たな開発者からの貢献を求めている。一部のプロジェクトでは、「Good First Issue(GFI)」と呼ばれるラベルを用いて、新規開発者向けの Issue を用意しオンボーディングを支援している。ただし、ラベル付けはプロジェクトメンテナの手作業で行われておりメンテナにとって負担となるため、多くのプロジェクトでは GFI ラベルは積極的に利用されていない。本研究の目的は、OSS プロジェクトの新規開発者向けの Issue を自動分類する機械学習モデルを構築することである。本論文では、ランダムフォレストを用いて分類モデルを構築した結果を述べる。通常 Issue 約 15 万件と GFI 約 1 万件を収集して 10 分割交差検証を行った結果、Precision が 0.91、Recall が 0.30 となった (RQ1)。また、重要度が高い特徴量を分析し、GFI の分類には投稿者のプロジェクト内での役割が重要であることが分かった (RQ2)。

Open Source Software (OSS) Development Projects are always looking for contributions from new developers to maintain the sustainability of the project. Some projects use a label called “Good First Issue (GFI)” to support onboarding by preparing issues for new developers. However, GFI labels are not initiatively used in many projects because labeling is done manually by project maintainers, which is burdensome for maintainers. The aim of this research is to construct a machine learning model to automatically classify issues for new developers in OSS projects. In this paper we describe the results of constructing a classification model using random forests method. We collected about 150,000 regular issues and about 10,000 GFIs, and conducted a 10-fold cross validation, resulting in a Precision of 0.91 and a Recall of 0.30 (RQ1). We also analyzed the feature value with high importance and found that the role of the contributor in a project is important for the classification of GFIs (RQ2).

## 1 はじめに

オープンソースソフトウェア (OSS) プロジェクトは、コミュニティからのバグ報告やバグ修正、機能拡張等の貢献によって成り立っている [1]。コミュニティ開発者らの多くはボランティアとして活動しており、自由に開発に参加、離脱することができるため、一般的な OSS プロジェクトでは、プロジェクトを持

続可能にするために常に新規開発者を求めている [2]。しかしながら、新規開発者はコミュニティに参加する際様々な障壁に直面し、プロジェクトから離脱する場合がある [7]。

新規開発者の直面する障壁の 1 つに課題選択の障壁がある [6]。多くの OSS プロジェクトでは貢献を行う際、開発者自身が課題票を作成するか、既知の課題の中から選択して取り組む。新規開発者にとっては、不慣れなプロジェクトの課題を理解し、対処するには多くの時間がかかる。

新規開発者がプロジェクトに参加する際の障壁を軽減するために、一部の OSS プロジェクトでは、新規開発者向けの課題に「Good First Issue(GFI)<sup>†1</sup>」と

Automatic Classification of Good First Issues for Onboarding to Open Source Software Projects.

Hyuga Horiguchi, 和歌山大学大学院システム工学研究科, Graduate School of Systems Engineering, Wakayama University.

Masao Ohira, 和歌山大学システム工学部, Faculty of Systems Engineering, Wakayama University.

コンピュータソフトウェア, Vol.39, No.4 (2022), pp.31–37. [研究論文 (レター)] 2022 年 2 月 11 日受付。

<sup>†1</sup> <https://docs.github.com/en/issues/using-labels-and-milestones-to-track-work/managing-labels>

呼ばれるラベルを付与している。実際に、GitHub 上で GFI ラベルを使用しているプロジェクトは過去 10 年間で増加しており [8], GFI の多くは新規開発者によって解決されている [5][4]。ただし、プロジェクトメンテナは課題を手動で選別しラベル付けを行う必要があり、メンテナにとって負担となる<sup>†2</sup>。

本研究の目的は、新規開発者向けの課題を分類する機械学習モデルを構築することである。GitHub 上に存在する GFI ラベルが付与された Issue (以降 GFI と呼ぶ) と通常の Issue を収集し、ランダムフォレストを使って機械学習を行う。構築したモデルを評価するために、以下のリサーチクエストに取り組む。

- **RQ1:** Issue 投稿時に含まれる情報のみを利用して GFI を分類できるか?
- **RQ2:** 分類に重要な Issue の特徴量は何か?

RQ1 では、分類モデルの性能を定量的に評価することを目的とする。メンテナの GFI ラベル付けの負担を軽減するため、分類モデルはメンテナが Issue の詳細を読みラベルを付与する前に、GFI かどうかを判定できる必要がある。そのため、モデルに入力する特徴量として利用できるのは、タイトルと本文、Issue 投稿者の情報のみであり、これらの特徴量を使って学習した際の分類モデルの性能を評価する。

RQ2 では、学習したモデルが何を基準に GFI と通常 Issue を分類したのか理解することを目的とする。ランダムフォレストでは、モデルに入力する特徴量が分類結果に与える相対的な重要度を計算することができる。重要度が高い特徴量は、GFI と通常 Issue を分ける特徴を理解するうえで役に立つ。

## 2 分類モデルの構築

### 2.1 入力する特徴量

本モデルは、メンテナが手動で行っている GFI のラベル付けを自動化することを目的としているため、メンテナが Issue の内容を確認する前の、Issue が投稿された直後の情報を使って GFI かどうかを判定できることが求められる。そこで、Issue 投稿時に必ず含まれる、タイトル、本文、投稿者の情報のみを使っ

て特徴量を考える。すなわち、Issue に対する他の開発者のコメントや、Issue に付与されるラベルの情報は特徴量として利用しない。モデルに入力する特徴量を表 1 に示す。

**タイトルと本文** タイトルと本文は、1 単語ずつに分割した後 Bag-of-Words としてベクトル化し、特徴量として利用する。また、タイトルと本文の文章の長さは Issue の内容が詳細に説明されているかどうかに関連するため、文章に含まれる単語数も特徴量として利用する。

**投稿者** プロジェクトオーナーやメンテナは、新規開発者を惹きつけるために、GFI になりえる Issue を積極的に投稿する可能性がある。そこで、Issue 投稿者がプロジェクトに対して持っている権限レベルをカテゴリ変数として数値化する。

### 2.2 タイトルと本文の前処理

GitHub 上の Issue のタイトルや本文は、Markdown 形式で記述することができる。また本文には、コード片や URL、バージョン情報などが含まれる。これらの情報は Issue の内容を理解するのに役立つため、正規表現を用いて検出し特定の単語に置換する。また自然言語に対しては、1 単語ごとのトークン化、小文字化、ストップワードと記号の除去、日付と数字の置換、レマタイズを行った。

### 2.3 分類モデル

分類モデルの機械学習アルゴリズムには、ランダムフォレストを利用する。モデルは Issue ごとの表 1 の特徴量を入力とし、GFI か通常 Issue かの 2 値分類を行う。ランダムフォレストでは、弱学習器として複数の決定木を学習する。個々の決定木は、元のデータセットをブートストラップサンプリングしたデータによって学習し、その際全特徴量のうちランダムサンプリングした特徴量のみを利用する。

決定木の個数と各決定木の特徴量の個数はハイパーパラメータであるため、本研究では [3] で推奨されている値を参考に、決定木の個数を 200、各決定木の特徴量の個数を  $\sqrt{F}$  とした。 $F$  は全特徴量の個数である。また本モデルの最終的な分類結果は、200 個の決

<sup>†2</sup> <https://github.blog/2020-01-22-how-we-built-good-first-issues/>

表 1 モデルに入力する特徴量

特徴量の名前	説明	
Title	BoW-n	タイトルの Bag-of-Words (n はタイトルの語彙数)
	Words	タイトルに含まれる単語数
Body	BoW-n	本文の Bag-of-Words (n は本文の語彙数)
	Words	本文に含まれる単語数
AuthorAssociation	Issue 投稿者が持つプロジェクトに対する権限レベル	

定木のうち過半数が GFI と判定した場合を GFI とし  
て分類する。

## 2.4 不均衡データの学習

3.1 節で示すように、GFI は通常 Issue と比べてサ  
ンプル数が非常に少なく不均衡なデータセットにな  
るため、クラスの重みづけを行う。決定木では、ノ  
ードを分割することで不純度が減少すると、その分割は  
分類に効果的であるとみなされる。そこで、不純度を  
計算する際に各クラスのサンプル数に反比例する重  
み  $w_i = \frac{N}{2N_i}$  ( $i$  はクラス、 $N$  は全サンプル数、 $N_i$  は  
クラス  $i$  のサンプル数) を掛けることで、サンプル数  
の少ない GFI の分類結果が不純度へ大きく影響する  
ようになり、結果として GFI の分類性能が向上する  
ことを期待できる。

## 3 実験方法

### 3.1 データセット

GFI ラベルを運用しているプロジェクトから  
GFI (正例) と通常 Issue (負例) を収集するため、以下  
の条件を満たす GitHub 上の Issue を全て取得し、そ  
れぞれの GFI が属するプロジェクトを特定した。

- 「good first issue」ラベルが付与されている
- 閉じられている (Close されている)

その結果、1 つ以上の GFI を保有している約 7 万件  
のプロジェクトを特定できた。これらすべてのプロ  
ジェクトから全ての Issue を収集するのは非常に時  
間がかかる上、GFI と通常 Issue のサンプル数の偏  
りを可能な限り小さくするため、GFI 数が 500 件以  
上の計 15 プロジェクトのみを対象とした。収集期間  
は各プロジェクト作成時から 2021 年 6 月 4 日まで  
である。

また、今回は GFI ラベルを付与していないプロジェ  
クトにも適用可能な、汎用性の高いモデルの構築を意  
図して、プロジェクト毎にはモデルを作らない。十分  
なデータが存在するプロジェクトに関しては、汎用モ  
デルよりも精度の高いモデルを個別に作成できる可  
能性があるため、今後詳細に調べる予定である。

### 3.2 評価

#### 3.2.1 RQ1: Issue 投稿時に含まれる情報のみ を利用して GFI を分類できるか?

データセットを 9:1 に分割し、9 割を使って 10 分  
割交差検証を行う。また、交差検証で学習したモデル  
の中から 1 つを使って残り 1 割のデータセットを分  
類し、混同行列と後の RQ2 で利用する特徴量の重要  
度を求める。評価指標には Precision, Recall, F1 を  
使用する。Precision は、モデルが GFI であると分類  
した Issue のうち、実際に GFI である Issue の割合  
を表す。Recall は、実際に GFI である Issue のうち、  
モデルが GFI であると分類した Issue の割合を表す。  
F1 は Precision と Recall の調和平均である。

#### 3.2.2 RQ2: 分類に重要な Issue の特徴量は 何か?

特徴量の重要度を用いて、分類結果に対する特徴量  
の相対的な重要性の評価を行う。重要度は、決定木の  
ノードをある特徴量を用いて分割した際、その分割に  
よって減少するノードの不純度と分割されるサンプ  
ル数の積で求められる。ここで不純度は、ノードに異  
なるクラスのサンプルがどの程度混在しているかを  
表す。したがって、分割によって不純度を大きく減少  
させる特徴量であるほど重要度が高くなり、また分割  
に影響するサンプル数が多いほど重要度が高くなる。  
重要度は決定木によって異なり、ランダムフォレスト

表 2 10 分割交差検証の結果

評価指標	平均
Precision	0.91
Recall	0.30
F1	0.45

表 3 評価用データの分類結果

		分類結果	
		GFI	通常 Issue
実際の	GFI	380	843
クラス	通常 Issue	38	14,867

の場合は決定木が複数あるため、平均値を用いる。

## 4 結果

### 4.1 RQ1: Issue 投稿時に含まれる情報のみを利用して GFI を分類できるか?

10 分割交差検証の結果を表 2 に示す。Precision の 0.91 と比べて Recall は 0.30 と低く、これはモデルが GFI であると分類したうち 9 割は実際に GFI であるが、7 割の GFI を見逃しているということを示している。また、10 個のモデルのうち 1 つを用いて評価用データを分類した結果、Precision は 0.91、Recall は 0.31、F1 は 0.46 で、混合行列は表 3 となった。

### 4.2 RQ2: 分類に重要な Issue の特徴量は何か?

特徴量の重要度を表 4 に示す。全特徴量の重要度の合計が 1 になるように正規化してある。また、テキストデータの Bag-of-Words の特徴量は、名前の「BoW-」の後に続く単語の重要度となっている。

最も重要度が高かったのは AuthorAssociation で、2 番目の Body Words に比べて 1.4 倍重要度が高かった。AuthorAssociation は、Issue 投稿者がプロジェクトに対してどのような権限を持っているかを示す特徴量で、Issue 投稿者のプロジェクト上での役割を意味する。プロジェクトのオーナーやメンテナは、プロジェクトの保守・運用のため新規開発者を惹きつけることに積極的であると考えられ、GFI になりえる Issue を通常の開発者よりも多く投稿している可能性

表 4 重要度上位 10 位までの特徴量

特徴量の名前	重要度
AuthorAssociation	0.0156
Body Words	0.0115
Title Words	0.0095
Body BoW-INLINE_CODE	0.0065
Body BoW-URL	0.0056
Body BoW-issue	0.0044
Body BoW-TIMESTAMP	0.0042
Body BoW-please	0.0039
Body BoW-VERSION	0.0037
Body BoW-add	0.0037

がある。

重要度が 2, 3 番目の本文とタイトルの単語数は、2.1 節で述べた通り、Issue の内容が詳細に説明されているかを表す特徴量である。また、Issue 本文に含まれるインラインコードや URL、バージョン情報の有無も重要度 10 位以内の特徴量に含まれており、これらは投稿された Issue の内容を深く理解するために重要な情報であるといえる。これらの情報が含まれていない Issue は、解決するための十分な情報が得られないため、新規開発者が取り組むべき Issue としては不向きであると考えられる。

## 5 議論

### 5.1 本モデルの有用性

4.1 節で示した通り、本モデルは Precision が 0.91 であるのに比べて Recall は 0.30 と非常に低い。Precision が高いということは、GFI と分類された Issue の中に通常 Issue が紛れ込むことは少なく、確実に新規開発者に適した Issue を分類することができることを意味する。しかしながら、Recall が低いということは、GFI になりえる Issue を見逃しやすいため、新規開発者がプロジェクトに貢献する機会が減る可能性がある。Tan ら [8] によれば、GFI ラベルはメンテナがサポート可能な Issue に対してラベル付けされるため、実際には多くの Issue を GFI に分類できたとしても、メンテナの手が回らず新規開発者はサポートが得られない可能性がある。したがって、多少 Recall を犠牲にしても Precision が高いほうが適している。

表 5 GFI, 通常 Issue 投稿者の AuthorAssociation の分布

AuthorAssociation	GFI[件] (割合)	通常 Issue[件] (割合)
プロジェクトのオーナー	<b>662 (0.05)</b>	<b>15 (0.00)</b>
プロジェクトを所有する組織のメンバー	3,662 (0.30)	21,873 (0.14)
プロジェクトの共同作業に招待された開発者	772 (0.06)	2,488 (0.02)
過去にこのプロジェクトにコミットしたことがある開発者	4,560 (0.37)	61,392 (0.41)
上記のどれにも当てはまらない開発者	2,727 (0.22)	65,227 (0.43)
合計	12,383	150,995

図 1 タイトルの単語数の分布

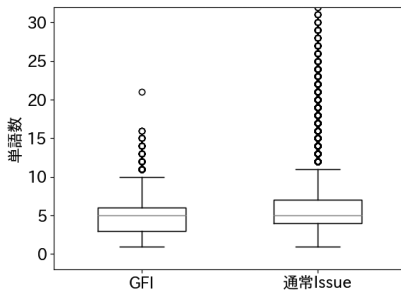
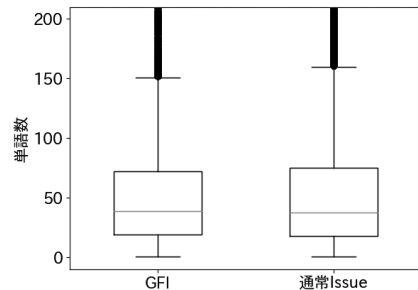


図 2 本文の単語数の分布



## 5.2 GFI の特徴

### 5.2.1 AuthorAssociation の分布

特徴量の重要度が 1 番高かった AuthorAssociation の分布を表 5 に示す。通常 Issue の合計件数 (150,995 件) は GFI の合計件数 (12,383 件) より圧倒的に多いにもかかわらず、プロジェクトオーナーが通常 Issue を投稿した件数はわずか 15 件で、GFI は 662 件投稿されていた。また、GFI の 41% がプロジェクトに対して何らかの権限を持っている開発者 (表 5 の上から 3 行分) から投稿されたのに対して、通常 Issue は 16% であった。したがって、プロジェクトのコアメンバーは、通常プロジェクトの課題よりも新規開発者向けの課題に対して多くの労力を割いており、本モデルを利用して Issue の中から新規開発者向けの課題を自動で抽出することで、メンテナの負担を軽減できる。

### 5.2.2 タイトルと本文の単語数の分布

特徴量の重要度が 2, 3 番目に高かった、タイトルと本文の単語数 (特徴量の名前は Title Words と Body Words) の箱ひげ図をそれぞれ図 1, 2 に示す。図 1 の箱ひげ図は、つぶれて見にくくなるのを防ぐために、

縦軸の最大値を 30 単語に制限した。同様の理由で、図 2 の縦軸の最大値を 200 単語に制限した。タイトルの単語数の中央値は、GFI は 5 単語で通常 Issue も 5 単語だった。また、本文の単語数の中央値は、GFI は 39 単語で、通常 Issue は 38 単語だった。GFI と通常 Issue で分布に大きな差がないにもかかわらず、特徴量の重要度が高くなった理由については、次の 2 つが考えられる。1 つは、タイトルや本文の単語数に加えて他の特徴量と組み合わせることで、分類性能が向上している場合である。もう 1 つは、外れ値を持つデータを分類する場合、分割時にノードに含まれる他クラスのデータがなく不純度が下がりやすいため、必然的に重要度が高くなる場合である。図 1, 2 を見ると外れ値を含むデータが多いことが分かるため、後者の場合を想定し、四分位範囲の 1.5 倍を超える単語数を持つデータを外れ値としてデータセットから除外し再度学習を行ったところ、Precision と Recall は変化しなかったが、タイトルと本文の単語数の重要度が上位 10 位から外れた。したがって、タイトルと本文の単語数は外れ値をとるデータを分類するためには有効であるが、GFI を特徴づける要素ではなかった。

### 5.3 制約

#### 5.3.1 GFI に関連する類似ラベル

本研究では、データ収集時に「beginner friendly」や「easy bug fix」といった「good first issue」に類似するラベルを含めなかった。これらの類似ラベルは「good first issue」と比べて使用頻度が低く [8]、特定のプロジェクトでしか使われないことが多いため、「good first issue」ラベルのみを収集した。

#### 5.3.2 収集した OSS プロジェクトの妥当性

データ収集の対象となった 15 プロジェクトのうち 2 件は、Issue の大半を GFI が占めており、通常 Issue はほとんど存在しなかった。これらのプロジェクトを目視したところ、一般に利用されることを想定したソフトウェア開発プロジェクトではなかったため、新規開発者が OSS プロジェクトに貢献する際の主要な動機であるスキルアップにはつながらない可能性が高く、多くの新規開発者にとっては GFI として不適當かもしれない。ただし、本研究では最初に 3.1 節で決めた条件にしたがってデータ収集を行い、恣意的にプロジェクトを除外することは避けた。

## 6 まとめ

本研究では、GFI を分類するランダムフォレストモデルを構築し、Precision は 0.91、Recall は 0.30 を達成した (RQ1)。また、重要度が高い特徴量について追加分析を行い、GFI は通常 Issue に比べてプロジェクトの保守や運用にかかわる開発者らによって投稿されていることが多いということが分かった (RQ2)。

つまりプロジェクトメンテナは、通常のプロジェクトの課題よりも新規開発者向けの課題に対して多くの労力を割いており、本モデルを利用して Issue の中から新規開発者向けの課題を自動で分類することは、メンテナの負担を軽減するために意義があるといえる。

現状のモデルは Recall が低いいため、GFI になりえる Issue を多く見逃している可能性がある。今後は、Recall が低い原因を分析し向上させることを目指す。

**謝辞** 本研究の一部は、文部科学省科学研究補助金 (基盤 (C): 22K11974) による助成を受けた。

### 参考文献

- [1] Crowston, K., Annabi, H., and Howison, J.: Defining Open Source Software Project Success, in *Proceedings of 24th the International Conference on Information Systems*, 2003, pp. 327–340.
- [2] Forte, A., and Lampe, C.: Defining, Understanding, and Supporting Open Collaboration: Lessons From the Literature, *American Behavioral Scientist*, Vol. 57, No. 5 (2013), pp. 535–547.
- [3] Hastie, T., Tibshirani, R., and Friedman, J.: *The Elements of Statistical Learning: Data Mining, Inference, and Prediction.*, Springer, second edition, 2017.
- [4] Horiguchi, H., Omori, I., and Ohira, M.: Onboarding to Open Source Projects with Good First Issues: A Preliminary Analysis, in *Proceedings of the 28th IEEE International Conference on Software Analysis, Evolution and Reengineering*, 2021, pp. 501–505.
- [5] Labuschagne, A., and Holmes, R.: Do Onboarding Programs Work?, in *Proceedings of 12th Working Conference on Mining Software Repositories*, 2015, pp. 381–385.
- [6] Steinmacher, I., Conte, T., and Gerosa, M. A.: Understanding and Supporting the Choice of an Appropriate Task to Start with in Open Source Software Communities, in *Proceedings of the 48th Hawaii International Conference on System Sciences*, 2015, pp. 5299–5308.
- [7] Steinmacher, I., Conte, T., Gerosa, M. A., and Redmiles, D.: Social Barriers Faced by Newcomers Placing Their First Contribution in Open Source Software Projects, in *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*, 2015, pp. 1379–1392.
- [8] Tan, X., Zhou, M., and Sun, Z.: A First Look at Good First Issues on GitHub, in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 398–409.



堀口日向

2020 年和歌山大学システム工学部システム工学科卒業、2022 年同大学院システム工学研究科システム工学専攻博士前期課程修了。修士 (工学)。ソフトウェア工学、特にオープンソースソフトウェアの新規開発者支援の研究に従事。



**大平 雅 雄**

1998年京都工芸繊維大学工芸学部電子情報工学科卒業，平成15年奈良先端科学技術大学院大学情報科学研究科博士後期課程修了。同大学情報科

学研究科助教を経て，2012年和歌山大学システム工学部講師（2014年より准教授）。博士（工学）。ソフトウェア工学，特にマイニングソフトウェアリポジトリの研究に従事。電子情報通信学会，情報処理学会，IEEE，ACM各会員。