

Temporal dynamics of Requirements Engineering from mobile app reviews

Vitor Mesaque Alves de Lima^{Corresp., 1}, Adailton Ferreira de Araújo², Ricardo Marcondes Marcacini^{1, 2}

¹ Faculty of Computing (FACOM), Federal University Of Mato Grosso do Sul (UFMS), Campo Grande, Mato Grosso do Sul, Brazil

² Institute of Mathematics and Computer Sciences (ICMC), University of São Paulo (USP), São Carlos, São Paulo, Brazil

Corresponding Author: Vitor Mesaque Alves de Lima
Email address: vitor.lima@ufms.br

Opinion mining for app reviews aims to analyze people's comments from app stores to support data-driven requirements engineering activities, such as bug report classification, new feature requests, and usage experience. However, due to a large amount of textual data, manually analyzing these comments is challenging, and machine-learning-based methods have been used to automate opinion mining. Although recent methods have obtained promising results for extracting and categorizing requirements from users' opinions, the main focus of existing studies is to help software engineers to explore historical user behavior regarding software requirements. Thus, existing models are used to support corrective maintenance from app reviews, while we argue that this valuable user knowledge can be used for preventive software maintenance. This paper introduces the temporal dynamics of requirements analysis to answer the following question: how to predict initial trends on defective requirements from users' opinions before negatively impacting the overall app's evaluation? We present the MAPP-Reviews (Monitoring App Reviews) method, which (i) extracts requirements with negative evaluation from app reviews, (ii) generates time series based on the frequency of negative evaluation, and (iii) trains predictive models to identify requirements with higher trends of negative evaluation. The experimental results from approximately 85,000 reviews show that opinions extracted from user reviews provide information about the future behavior of an app requirement, thereby allowing software engineers to anticipate the identification of requirements that may affect the future app's ratings.

Temporal dynamics of Requirements Engineering from mobile app reviews

Vitor Mesaque Alves de Lima¹, Adailton Ferreira de Araújo², and Ricardo Marcondes Marcacini^{1,2}

¹Faculty of Computing (FACOM), Federal University Of Mato Grosso do Sul (UFMS), Campo Grande, Mato Grosso do Sul, Brazil

²Institute of Mathematics and Computer Sciences (ICMC), University of São Paulo (USP), São Carlos, São Paulo, Brazil

Corresponding author:

Vitor Mesaque Alves de Lima¹

Email address: vitor.lima@ufms.br

ABSTRACT

Opinion mining for app reviews aims to analyze people's comments from app stores to support data-driven requirements engineering activities, such as bug report classification, new feature requests, and usage experience. However, due to a large amount of textual data, manually analyzing these comments is challenging, and machine-learning-based methods have been used to automate opinion mining. Although recent methods have obtained promising results for extracting and categorizing requirements from users' opinions, the main focus of existing studies is to help software engineers to explore historical user behavior regarding software requirements. Thus, existing models are used to support corrective maintenance from app reviews, while we argue that this valuable user knowledge can be used for preventive software maintenance. This paper introduces the temporal dynamics of requirements analysis to answer the following question: how to predict initial trends on defective requirements from users' opinions before negatively impacting the overall app's evaluation? We present the MAPP-Reviews (Monitoring App Reviews) method, which (i) extracts requirements with negative evaluation from app reviews, (ii) generates time series based on the frequency of negative evaluation, and (iii) trains predictive models to identify requirements with higher trends of negative evaluation. The experimental results from approximately 85,000 reviews show that opinions extracted from user reviews provide information about the future behavior of an app requirement, thereby allowing software engineers to anticipate the identification of requirements that may affect the future app's ratings.

INTRODUCTION

Opinions extracted from app reviews provide a wide range of user feedback to support requirements engineering activities, such as bug report classification, new feature requests, and usage experience (Dabrowski et al., 2020; Martin et al., 2016; AlSubaih et al., 2019; Araujo and Marcacini, 2021). However, manually analyzing a reviews dataset to extract useful knowledge from the opinions is challenging because of the large amount of data and the high frequency of new reviews published by users (Johanssen et al., 2019; Martin et al., 2016). To deal with these challenges, opinion mining has been increasingly used for computational analysis of the people's opinions from free texts (B. Liu, 2012). In the context of app reviews, opinion mining allows extracting excerpts from comments and mapping them to software requirements, as well as classifying the positive, negative or neutral polarity of these requirements according to the users' experience (Dabrowski et al., 2020).

One of the main challenges for software quality maintenance is identifying emerging issues, e.g., bugs, in a timely manner (April and Abran, 2012). These issues can generate huge losses, as users can fail to perform important tasks or generate dissatisfaction that leads the user to uninstall the app. A recent survey showed that 78.3% of developers consider removing unnecessary and defective requirements to be equally or more important than adding new requirements (Nayebi, Kuznetsov, et al., 2018). According to Lientz and Swanson (1980), maintenance activities are categorized into four classes: i) adaptive - changes

in the software environment; ii) perfective - new user requirements; iii) corrective - fixing errors; and iv) preventive - prevent problems in the future. The authors showed that around 21% of the maintenance effort was on the last two types (Bennett and Rajlich, 2000). Specifically, in the context of mobile apps Mcilroy, Ali, and Hassan (2016) found that rationale for the update most frequently communicated task in app stores is bug fixing which occurs in 63% of the updates. Thus, approaches that automate the analysis of potentially defective software requirements from app reviews are important to make strategic updates, as well as prioritization and planning of new releases (Licorish, Savarimuthu, and Keertipati, 2017). In addition, the app stores offer a more dynamic way of distributing the software directly to users, with shorter release times than traditional software systems, i.e., continuous update releases are performed every few weeks or even days (Nayebi, Adams, and Ruhe, 2016). Therefore, app reviews provide quick feedback from the crowd about software misbehavior that may not necessarily be reproducible during regular development/testing activities, e.g., device combinations, screen sizes, operating systems and network conditions (Palomba, Linares-Vásquez, Bavota, Oliveto, Penta, et al., 2018). This continuous crowd feedback can be used by developers in the development and preventive maintenance process.

Using an opinion mining approach, we argue that software engineers can investigate bugs and misbehavior early when an app receives negative reviews. Opinion mining techniques can organize reviews based on the identified software requirements and their associated user's sentiment (Dabrowski et al., 2020). Consequently, developers can examine negative reviews about a specific feature to understand the user's concerns about a defective requirement and potentially fix it more quickly, i.e., before impacting many users and negatively affecting the app's ratings.

Different strategies have recently been proposed to discover these emerging issues (Zhao et al., 2020), such as issues categorization (Tudor and Walter, 2006; Jacob and Harrison, 2013; Galvis Carreño and Winbladh, 2013; Pagano and W. Maalej, 2013; Mcilroy, Ali, Khalid, et al., 2016; Khalid et al., 2015; Panichella et al., 2015; Panichella et al., 2016), sentiment analysis of the software requirements to identify certain levels of dissatisfaction (Gao, Zeng, Wen, et al., 2020), and analyze the degree of utility of a requirement (Guzman and Walid Maalej, 2014). These approaches are concerned only with past reviews and acting in a corrective way, i.e., these approaches do not have preventive strategies to anticipate problems that can become frequent and impact more users in the coming days or weeks. Analyzing the temporal dynamics of a requirement from app reviews provides information about a requirement's future behavior. In this sense, we raise the following research question: how do we predict initial trends on defective requirements from users' opinions before negatively impacting the overall app's evaluation?

In this paper, we present the MAPP-Reviews (Monitoring App Reviews) method. MAPP-Reviews explores the temporal dynamics of software requirements extracted from app reviews. First, we collect, pre-process and extract software requirements from large review datasets. Then, the software requirements associated with negative reviews are organized into groups according to their content similarity by using clustering technique. The temporal dynamics of each requirement group is modeled using a time series, which indicates the time frequency of a software requirement from negative reviews. Finally, we train predictive models on historical time series to forecast future points. Forecasting is interpreted as signals to identify which requirements may negatively impact the app in the future, e.g., identify signs of app misbehavior before impacting many users and prevent the low app ratings. Our main contributions are briefly summarized below:

- Although there are promising methods for extracting candidate software requirements from application reviews, such methods do not consider that users describe the same software requirement in different ways with non-technical and informal language. Our MAPP-Reviews method introduces software requirements clustering to standardize different software requirement writing variations. In this case, we explore contextual word embeddings for software requirements representation, which have recently been proposed to support natural language processing. When considering the clustering structure, we can more accurately quantify the number of negative user mentions of a software requirement over time.
- We present a method to generate the temporal dynamics of negative ratings of a software requirements cluster by using time series. Our method uses equal-interval segmentation to calculate the frequency of software requirements mentions in each time interval. Thus, a time series is obtained and used to analyze and visualize the temporal dynamics of the cluster, where we are especially interested in intervals where sudden changes happen.

- Time series forecasting is useful to identify in advance an upward trend of negative reviews for a given software requirement. However, most existing forecasting models do not consider domain-specific information that affects user behavior, such as holidays, new app releases and updates, marketing campaigns, and other external events. In the MAPP-Reviews method, we investigate the incorporation of software domain-specific information through trend changepoints. We explore both automatic and manual changepoint estimation.

We carried out an experimental evaluation involving approximately 85,000 reviews over 2.5 years for three food delivery apps. The experimental results show that it is possible to find significant points in the time series that can provide information about the future behavior of the requirement through app reviews. Our method can provide important information to software engineers regarding software development and maintenance. Moreover, software engineers can act preventively through the proposed MAPP-Reviews approach and reduce the impacts of a defective requirement.

This paper is structured as follows. Section “Background and Related Work” presents the literature review and related work about mining user opinions to support requirement engineering and emerging issue detection. In “MAPP-Reviews method” section, we present the architecture of the proposed method. We present the main results in “Results” section. Thereafter, we evaluate and discuss the main findings of the research in “Discussion” section. Finally, in “Conclusions” section, we present the final considerations and future work.

BACKGROUND AND RELATED WORK

The opinion mining of app reviews can involve several steps, such as software requirements organization from reviews (Araujo and Marcacini, 2021), grouping similar apps using textual features (Al-Subaihin et al., 2016; Harman, Jia, and Yuanyuan Zhang, 2012), reviews classification in categories of interest to developers (e.g., Bug and New Features) (Araujo, Golo, et al., 2020), sentiment analysis of the users’ opinion about the requirements (Dragoni, Federici, and Rexha, 2019; Malik, Shakshuki, and Yoo, 2020), and the prediction of the review utility score (Ying Zhang and Lin, 2018). The requirements extraction has an essential role in these steps since the failure in this task directly affects the performance of the other steps.

Dabrowski et al. (2020) evaluated the performance of the three state-of-the-art requirements extraction approaches: SAFE (Johann, Stanik, Walid Maalej, et al., 2017), ReUS (Dragoni, Federici, and Rexha, 2019) and GuMa (Guzman and Walid Maalej, 2014). These approaches explore rule-based information extraction from linguistic features. GuMa (Guzman and Walid Maalej, 2014) used a co-location algorithm, thereby identifying expressions of two or more words that correspond to a conventional way of referring to things. SAFE (Johann, Stanik, Walid Maalej, et al., 2017) and ReUS (Dragoni, Federici, and Rexha, 2019) defined linguistic rules based on grammatical classes and semantic dependence. The experimental evaluation of (Dabrowski et al., 2020) revealed that the low accuracy presented by the rule-based approaches could hinder its use in practice.

Araujo and Marcacini (2021) proposed RE-BERT (Requirements Engineering using Bidirectional Encoder Representations from Transformers) method for software requirements extraction from reviews based on Local Context Word Embeddings (i.e. deep neural language model). RE-BERT models the requirements extraction as a token classification task from deep neural networks. To solve some limitations of rule-based approaches, RE-BERT allows the generation of word embeddings for reviews according to the context of the sentence in which the software requirement occurs. Moreover, RE-BERT explores a multi-domain training strategy to enable software requirements extraction from app reviews of new domains without labeled data.

After extracting requirements from app reviews, there is a step to identify more relevant requirements and organize them into groups of similar requirements. Traditionally, requirements obtained from user interviews are prioritized with manual analysis techniques, such as the MoSCoW (Tudor and Walter, 2006) method that categorizes each requirement into groups, and applies the AHP (Analytical Hierarchy Process) decision-making (Saaty, 1980). These techniques are not suitable for prioritizing large numbers of software requirements because they require domain experts to categorize each requirement. Therefore, recent studies have applied data mining approaches and statistical techniques (Pagano and W. Maalej, 2013).

The statistical techniques have been used to find issues such as to examine how app features predict an app’s popularity (M. Chen and X. Liu, 2011), to analyze the correlations between the textual size of

the reviews and users' dissatisfaction (Vasa et al., 2012), lower rating and negative sentiments (Hoon et al., 2012), correlations between the rating assigned by users and the number of app downloads (Harman, Jia, and Yuanyuan Zhang, 2012), to the word usage patterns in reviews (Gómez et al., 2015; Licorish, Savarimuthu, and Keertipati, 2017), to detect traceability links between app reviews and code changes addressing them (Palomba, Linares-Vásquez, Bavota, Oliveto, Penta, et al., 2018), and explore the feature lifecycles in app stores (Sarro et al., 2015). There also exists some work focus on defining taxonomies of reviews to assist mobile app developers with planning maintenance and evolution activities (Di Sorbo et al., 2016; Ciurumelea et al., 2017; Nayebi, Kuznetsov, et al., 2018). In addition to user reviews, previous works (Guzman, Alkadhi, and Seyff, 2016; Guzman, Alkadhi, and Seyff, 2017; Nayebi, Cho, and Ruhe, 2018) explored how a dataset of tweets can provide complementary information to support mobile app development.

From a labeling perspective, previous works classified and grouped software reviews into classes and categories (Jacob and Harrison, 2013; Galvis Carreño and Winbladh, 2013; Pagano and W. Maalej, 2013; Mcilroy, Ali, Khalid, et al., 2016; Khalid et al., 2015; N. Chen et al., 2014; Gómez et al., 2015; Gu and Kim, 2015; Walid Maalej and Nabil, 2015; Villarroel et al., 2016; Nayebi, Marbouti, et al., 2017), such as feature requests, requests for improvements, requests for bug fixes, and usage experience. Noei, F. Zhang, and Zou (2021) used topic modeling to determine the key topics of user reviews for different app categories.

Regarding analyzing emerging issues from app reviews, existing studies are usually based on topic modeling or clustering techniques. For example, LDA (Latent Dirichlet Allocation) (Blei, Ng, and Jordan, 2003), DIVER (iDentifying emerging app Issues Via usER feedback) (Gao, Zheng, et al., 2019) and IDEA (Gao, Zeng, Lyu, et al., 2018) approaches were used for app reviews. The LDA approach is a topic modeling method used to determine patterns of textual topics, i.e., to capture the pattern in a document that produces a topic. LDA is a probabilistic distribution algorithm for assigning topics to documents. A topic is a probabilistic distribution over words, and each document represents a mixture of latent topics (Guzman and Walid Maalej, 2014). In the context of mining user opinions in app reviews, especially to detect emerging issues, the documents in the LDA are app reviews, and the extracted topics are used to detect emerging issues. The IDEA approach improves LDA by considering topic distributions in a context window when detecting emerging topics by tracking topic variations over versions (Gao, Zeng, Wen, et al., 2020). In addition, the IDEA approach implements an automatic topic interpretation method to label each topic with the most representative sentences and phrases (Gao, Zeng, Wen, et al., 2020). In the same direction, the DIVER approach was proposed to detect emerging app issues, but mainly in beta test periods (Gao, Zheng, et al., 2019). The IDEA, DIVER and LDA approaches have not been considered sentiment of user reviews. Recently, the MERIT (iMproved EmeRging Issue deTectioN) (Gao, Zeng, Wen, et al., 2020) approach was proposed and explore word embedding techniques to prioritize phrases/sentences of each positive and negative topic. Phong et al. (2015) and Vu et al. (2016) grouped the keywords and phrases using clustering algorithms and then determine and monitor over time the emergent clusters based on the occurrence frequencies of the keywords and phrases in each cluster. Palomba, Linares-Vásquez, Bavota, Oliveto, Di Penta, et al. (2015) proposes an approach to tracking informative user reviews of source code changes and to monitor the extent to which developers addressing user reviews. These approaches are descriptive models, i.e., they analyze historical data to interpret and understand the behavior of past reviews. In our paper, we are interested in predictive models that aim to anticipate the growth of negative reviews that can impact the app's evaluation.

In short, app reviews formed the basis for many studies and decisions ranging from feature extraction to release planning of mobile apps. However, previous related works do not explore the temporal dynamics with a predictive model of requirements in reviews, as shown in Table 1. Related works that incorporate temporal dynamics cover only descriptive models. In addition, existing studies focus on only a few steps of the opinion mining process from app reviews, which hinders its use in real-world applications. Our proposal instantiates a complete opinion mining process and incorporates temporal dynamics of software requirements extracted from app reviews into forecasting models to address these drawbacks.

THE MAPP-REVIEWS METHOD

In order to analyze the temporal dynamics of software requirements, we present the MAPP-Reviews approach with five stages, as shown in Figure 1. First, we collect mobile app reviews in app stores through a web crawler. Second, we group the similar extracted requirements by using clustering methods. Third,

Table 1. Overview of related works.

Reference	Data Representation	Pre-processing and Extraction of Requirements	Requirements/Topics Clustering and Labeling	Temporal Dynamics
(Araujo and Marcacini, 2021)	Word embeddings.	Token Classification.	No.	No.
(Gao, Zeng, Wen, et al., 2020)	Word embeddings.	Rule-based and Topic modeling.	Yes. It combines word embeddings with topic distributions as the semantic representations of words.	Yes. Descriptive Model.
(Malik, Shakshuki, and Yoo, 2020)	Bag-of-words.	Rule-based.	No.	No.
(Gao, Zheng, et al., 2019)	Vector space.	Rule-based and Topic modeling.	Yes. Anomaly Clustering Algorithm.	Yes. Descriptive model.
(Dragoni, Federici, and Rexha, 2019)	Dependency tree.	Rule-based.	No.	No.
(Gao, Zeng, Lyu, et al., 2018)	Probability vector.	Rule-based and Topic modeling.	Yes. AOLD - Adaptively Online Latent Dirichlet Allocation. The topic labeling method considers the semantic similarity between the candidates and the topics.	Yes. Descriptive model.
(Johann, Stanik, Walid Maalej, et al., 2017)	Keywords.	Rule-based.	No.	No.
(Vu et al., 2016)	Word embeddings.	Pre-defined.	Yes. Soft Clustering algorithm that uses vector representation of words from Word2vec.	Yes. Descriptive model.
(Villarroel et al., 2016)	Bag-of-words.	Rule-based.	Yes. DBSCAN clustering algorithm. Each cluster has a label composed of the five most frequent terms.	No.
(Gu and Kim, 2015)	Semantic Dependence Graph.	Rule-based.	Yes. Clustering aspect-opinion pairs with the same aspects.	Yes. Descriptive model.
(Phong et al., 2015)	Vector space.	Rule-based.	Yes. Word2vec and K-means.	Yes. Descriptive model.
(Guzman and Walid Maalej, 2014)	Keywords.	Rule-based and Topic modeling.	Yes. LDA approach.	No.
(N. Chen et al., 2014)	Bag-of-words.	Topic modeling.	Yes. LDA and ASUM approach with labeling.	Yes. Descriptive model.
(Iacob and Harrison, 2013)	Keywords.	Rule-based and Topic modeling.	Yes. LDA approach.	No.
(Galvis Carreño and Winbladh, 2013)	Bag-of-words.	Topic modeling.	Yes. Aspect and Sentiment Unification Model (ASUM) approach.	No.
(Harman, Jia, and Yuanyuan Zhang, 2012)	Keywords.	Pre-defined.	Yes. Greedy-based clustering algorithm.	No.
(Palomba, Linares-Vásquez, Bavota, Oliveto, Penta, et al., 2018)	Bag-of-words.	Topic-modeling.	Yes. AR-Miner approach with labeling.	No.

the most relevant clusters are identified to generate time series from negative reviews. Finally, we train the predictive model from time series to forecast software requirements involved with negative reviews, which will potentially impact the app's rating.

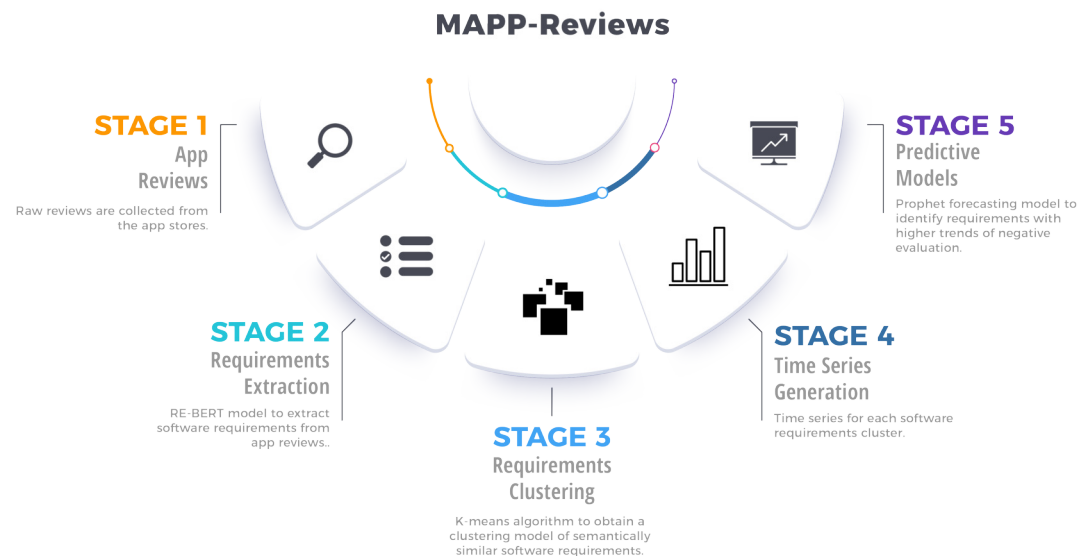


Figure 1. Overview of the proposed method for analyzing temporal dynamics of requirements engineering from mobile app reviews.

App Reviews

The app stores provide the textual content of the reviews, the publication date, and the rating stars of user-reported reviews. In the first stage of MAPP-Reviews, raw reviews are collected from the app stores using a web crawler tool through a RESTful API. At this stage, there is no pre-processing in the textual content of reviews. Data is organized in the appropriate data structure and automatically batched to be processed by the requirements extraction stage of MAPP-Reviews. In the experimental evaluation presented in this article, we used reviews collected from three food delivery apps: Uber Eats, Foodpanda, and Zomato.

Requirements Extraction

This section describes stages 2 of the MAPP-Reviews method, where there is the software requirements extraction from app reviews and text pre-processing using contextual word embeddings.

MAPP-Reviews uses the pre-trained RE-BERT (Araujo and Marcacini, 2021) model to extract software requirements from app reviews. RE-BERT is an extractor developed from our previous research. We trained the RE-BERT model using a labeled reviews dataset generated with a manual annotation process, as described by Dabrowski et al. (2020). The reviews are from 8 apps of different categories as showed in Table 2. RE-BERT uses a cross-domain training strategy, where the model was trained in 7 apps and tested in one unknown app for the test step. RE-BERT software requirements extraction performance was compared to SAFE (Johann, Stanik, Walid Maalej, et al., 2017), ReUS (Dragoni, Federici, and Rexha, 2019) and GuMa (Guzman and Walid Maalej, 2014). Since RE-BERT uses pre-trained models for semantic representation of texts, the extraction performance is significantly superior to the rule-based methods. Given this scenario, we selected RE-BERT for the requirement extraction stage. Figure 2 shows an example of review and extracted software requirements. In the raw review “I am ordering with delivery but it is automatically placing order with pick-up”, four software requirements were extracted (“ordering”, “delivery”, “placing order”, and “pick-up”). Note that “placing order” and “ordering” are the same requirement in practice. In the clustering step of the MAPP-Reviews method, these requirements are grouped in the same cluster, as they refer to the same feature.

RE-BERT returns the probability that each token (e.g. word) is a software requirement. Consecutive tokens in a sentence are concatenated to obtain software requirements expressions composed of two or

Table 2. Statistics about the datasets from 8 apps of different categories used to train the RE-BERT model.

	eBay	Evernote	Facebook	Netflix	Photo editor	Spotify	Twitter	WhatsApp
Reviews	1,962	4,832	8,293	14,310	7,690	14,487	63,628	248,641
Category	Shopping	Productivity	Social	Entertainment	Photography	Music and Audio	Social	Communication



Figure 2. Example of a review and extracted requirements.

more tokens. We filter reviews that are more associated with negative comments through user feedback. Consider that the user gives a star rating when submitting a review for an app. Generally, the star rating ranges from 1 to 5. This rating can be considered as the level of user satisfaction. In particular, we are interested in defective software requirements, and only reviews with 1 or 2 rating stars were considered. Thus, we use RE-BERT to extract only software requirements mentioned in reviews that may involve complaints, bad usage experience, or malfunction of app features.

RE-BERT extracts software requirements directly from the document reviews and we have to deal with the drawback that the same requirement can be written in different ways by users. Thus, we propose a software requirement semantic clustering, in which different writing variations of the same requirement must be standardized. However, the clustering step requires that the texts be pre-processed and structured in a format that allows the calculation of similarity measures between requirements.

We represent each software requirement through contextual word embedding. Word embeddings are vector representations for textual data in an embedding space, where we can compare two texts semantically using similarity measures. Different models of word embeddings have been proposed, such as Word2vec (Mikolov et al., 2013), Glove (Pennington, Socher, and Manning, 2014), FastText (Bojanowski et al., 2017) and BERT (Devlin et al., 2018). We use the BERT Sentence-Transformers model (Reimers and Gurevych, 2019) to maintain a neural network architecture similar to RE-BERT. BERT is a contextual neural language model, where for a given sequence of tokens, we can learn a word embedding representation for a token. Word embeddings can calculate the semantic proximity between tokens and entire sentences, and the embeddings can be used as input to train the classifier. BERT-based models are promising to learn contextual word embeddings from long-term dependencies between tokens in sentences and sentences (Araujo and Marcacini, 2021). However, we highlight that a local context more impacts the extraction of software requirements from reviews, i.e., tokens closer to those of software requirements are more significant (Araujo and Marcacini, 2021). Therefore, RE-BERT explores local contexts to identify relevant candidates for software requirements. Formally, let $E = \{r_1, r_2, \dots, r_n\}$ be a set of n extracted software requirements, where $r_i = (t_1, \dots, t_k)$ are a sequence of k tokens of the requirement r_i . BERT explores a masked language modeling procedure, i.e., BERT model first generates a corrupted \hat{x} version of the sequence, where approximately 15% of the words are randomly selected to be replaced by

a special token called [MASK] (Araujo and Marcacini, 2021). One of the training objectives is the noisy reconstruction defined in Equation 1,

$$p(\bar{r}|\hat{r}) = \sum_{j=1}^k m_j \frac{\exp(\mathbf{h}_{c_j}^\top \mathbf{w}_{t_j})}{\sum_{t'} \exp(\mathbf{h}_{c_j}^\top \mathbf{w}_{t'})} \quad (1)$$

where \hat{r} is a corrupted token sequence of requirement r , \bar{r} is the masked tokens, m_t is equal to 1 when t_j is masked and 0 otherwise. The c_t represents context information for the token t_j , usually the neighboring tokens. We extract token embeddings from the pre-trained BERT model, where \mathbf{h}_{c_j} is a context embedding and \mathbf{w}_{t_j} is a word embedding of the token t_j . The term $\sum_{t'} \exp(\mathbf{h}_{c_j}^\top \mathbf{w}_{t'})$ is a normalization factor using all tokens t' from a context c . BERT uses the Transformer deep neural network to solve $p(\bar{r}|\hat{r})$ of the Equation 1. Figure 3 illustrates a set of software requirements in a two-dimensional space obtained from contextual word embeddings. Note that the vector space of embeddings preserves the proximity of similar requirements, but written in different ways by users such as “search items”, “find items”, “handles my searches” and “find special items”.

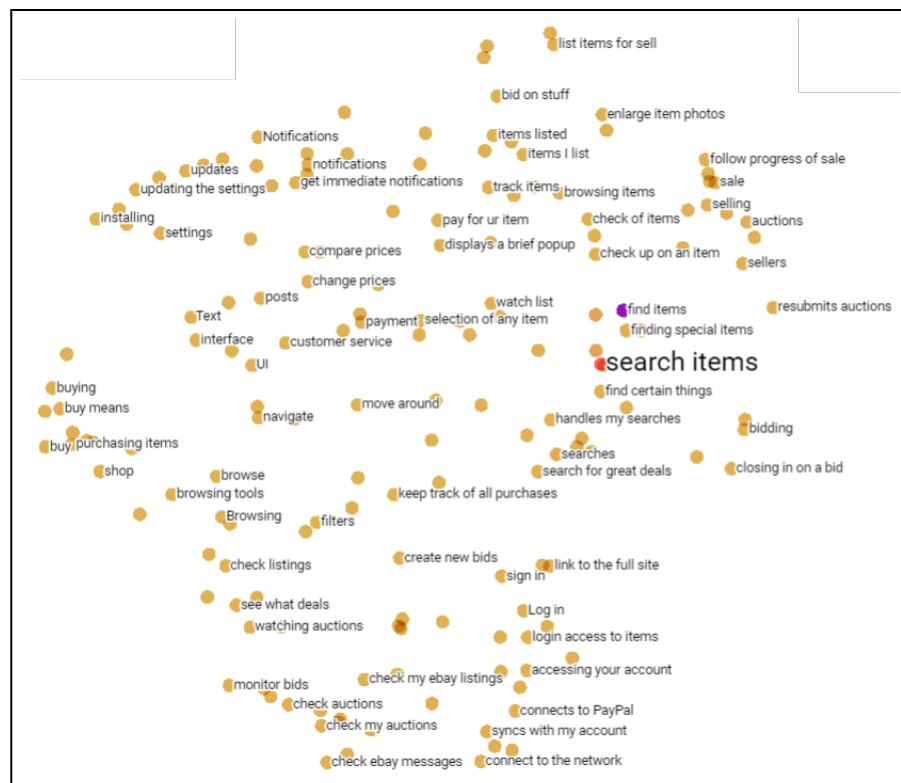


Figure 3. Set of software requirements in a two-dimensional space obtained from contextual word embeddings.

Requirements Clustering

After mapping the software requirements into word embeddings, MAPP-Reviews uses the k-means algorithm (MacQueen et al., 1967) to obtain a clustering model of semantically similar software requirements.

Formally, let $R = \{r_1, r_2, \dots, r_n\}$ a set of extracted software requirements, where each requirement r is a m -dimensional real vector from an word embedding space. The k-means clustering aims to partition the n requirements into k ($2 \leq k \leq n$) clusters $C = \{C_1, C_2, \dots, C_k\}$, thereby minimizing the within-cluster sum of squares as defined in Equation 2, where μ_i is the mean vector of all requirements in C_i .

$$\sum_{C_i \in C} \sum_{r \in C_i} \|r - \mu_i\|^2 \quad (2)$$

We observe that not all software requirements cluster represents a functional requirement in practice. Then, we evaluated the clustering model using a statistical measure called silhouette (Rousseeuw, 1987) to discard clusters with many different terms and irrelevant requirements. The silhouette value of a data instance is a measure of how similar a software requirement is to its own cluster compared to other clusters. The silhouette measure ranges from -1 to $+1$, where values close to $+1$ indicate that the requirement is well allocated to its own cluster (Vendramin, Campello, and Hruschka, 2010). Finally, we use the requirements with higher silhouette values to support the cluster labeling, i.e., to determine the software requirement's cluster name. For example, Table 3 shows the software requirement cluster "Payment" and some tokens allocated in the cluster with their respective silhouette values.

Table 3. Example of software requirement cluster "Payment" and some tokens allocated in the cluster with their respective silhouette values.

Cluster Label	Tokens with Silhouette (s)
Payment	"payment getting" ($s = 0.2618$), "payment get" ($s = 0.2547$), "getting payment" ($s = 0.2530$), "take payment" ($s = 0.2504$), "payment taking" ($s = 0.2471$), "payment" ($s = 0.2401$)

To calculate the silhouette measure, let $r_i \in C_i$ a requirement r_i in the cluster C_i . Equation 3 compute the mean distance between r_i and all other software requirements in the same cluster, where $d(r_i, r_j)$ is the distance between requirements r_i and r_j in the cluster C_i . In the equation, the expression $\frac{1}{|C_i|-1}$ means the distance $d(r_i, r_i)$ is not added to the sum. A smaller value of the silhouette measure $a(i)$ indicates that the requirement i is far from neighboring clusters and better assigned to its cluster.

$$a(r_i) = \frac{1}{|C_i|-1} \sum_{r_j \in C_i, r_i \neq r_j} d(r_i, r_j) \quad (3)$$

Analogously, the mean distance from requirement r_i to another cluster C_k is the mean distance from r_i to all requirements in C_k , where $C_k \neq C_i$. For each requirement $r_i \in C_i$, Equation 4 defines the minimum mean distance of r_i for all requirements in any other cluster, of which r_i is not a member. The cluster with this minimum mean distance is the neighbor cluster of r_i . So this is the next best-assigned cluster for the r_i requirement. The silhouette (value) of the software requirement r_i is defined by Equation 5.

$$b(r_i) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{r_j \in C_k} d(r_i, r_j) \quad (4)$$

$$s(r_i) = \frac{b(r_i) - a(r_i)}{\max\{a(r_i), b(r_i)\}}, \text{ if } |C_i| > 1 \quad (5)$$

At this point in the MAPP-Reviews method, we have software requirements pre-processed and represented through contextual word embeddings, as well as an organization of software requirements into k clusters. In addition, each cluster has a representative text (cluster label) obtained according to the requirements with higher silhouette values.

Figure 4 shows a two-dimensional projection of clustered software requirements from approximately 85,000 food delivery app reviews, which were used in the experimental evaluation of this work. High-density regions represent clusters of similar requirements that must be mapped to the same software requirement during the analysis of temporal dynamics. In the next section, techniques for generating the time series from software requirements clusters are presented, as well as the predictive models to infer future trends.

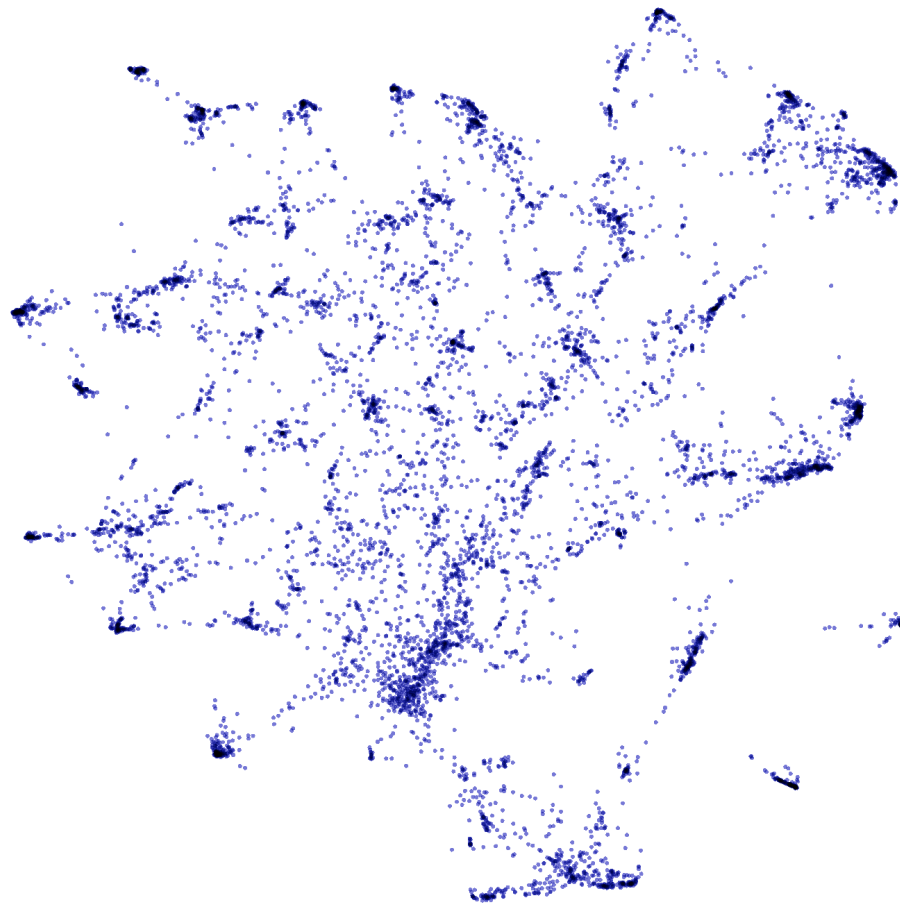


Figure 4. Two-dimensional projection of clustered software requirements from approximately 85,000 food delivery app reviews.

Time Series Generation

Time series can be described as an ordered sequence of observations (Chatfield and Xing, 2019). A time series of size s is defined as $X = (x_1, x_2, \dots, x_s)$ in which $x_t \in \mathbb{R}$ represents an observation at time t .

MAPP-Reviews generates time series for each software requirements cluster, where the observations represent how many times each requirement occurred in a period. Consequently, we know how many times a specific requirement was mentioned in the app reviews for each period. Each series models the temporal dynamics of a software requirement, i.e., the temporal evolution considering occurrences in negative reviews.

Some software requirements are naturally more frequent than others, as well as the tokens used to describe these requirements. For the time series analysis to be compared uniformly, we generate a normalized series for each requirement. Each observation in the time series is normalized according to Equation 6,

$$x_{normalized} = \frac{x}{z_p} \quad (6)$$

where $x_{normalized}$ is the result of the normalization, where x is the frequency of cluster (time series observation) C in the period p , z_p is the total frequency of the period.

Figure 5 shows an example of one of the generated time series for a software requirement. The time dynamics represented in the time series indicate the behavior of the software requirement concerning negative reviews. Note that in some periods there are large increases in the mention of the requirement,

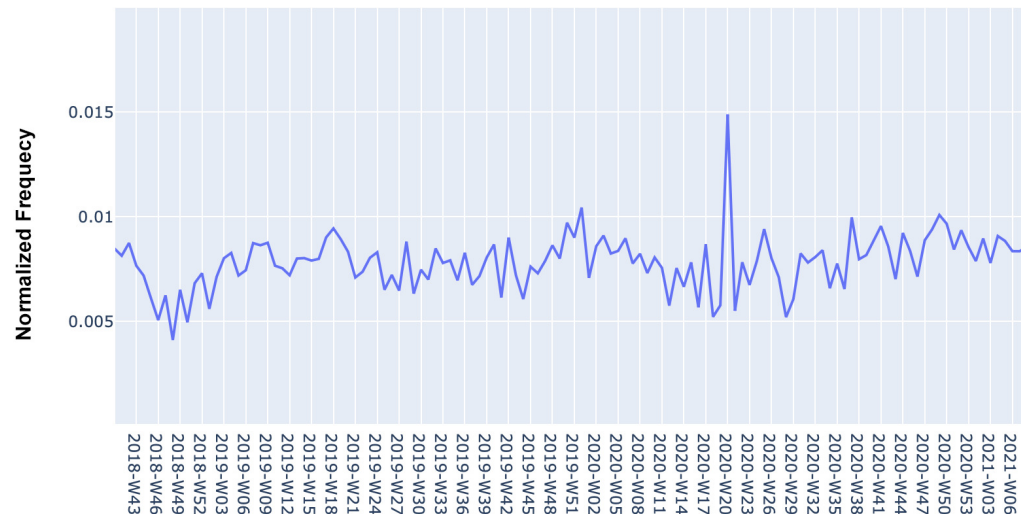


Figure 5. Time series with the normalized frequency of “Arriving time” requirement from Zomato App in negative reviews.

thereby indicating that users have negatively evaluated the app for that requirement. Predicting the occurrence of these periods for software maintenance, aiming to minimize the number of future negative reviews is the objective of the MAPP-Reviews predictive model discussed in the next section.

Predictive Models

Predictive models for time series are very useful to support an organization in its planning and decision-making. Such models explore past observations to estimate observations in future horizons, given a confidence interval. In our MAPP-Reviews method, we aim to detect the negative reviews of a software requirement that are starting to happen and make a forecast to see if they will become serious in the subsequent periods, i.e., high frequency in negative reviews. The general idea is to use p points from the time series to estimate the next $p + h$ points, where h is the prediction horizon.

MAPP-Reviews uses the Prophet Forecasting Model (Taylor and Letham, 2018). Prophet is a model from Facebook researchers for forecasting time series data considering non-linear trends at different time intervals, such as yearly, weekly, and daily seasonality. We chose the Prophet model for the MAPP-Reviews method due to the ability to incorporate domain knowledge into the predictive model. The Prophet model consists of three main components, as defined in Equation 7,

$$y(t) = g(t) + s(t) + h(t) + t\epsilon \quad (7)$$

where $g(t)$ represents the trend, $s(t)$ represents the time series seasonality, $h(t)$ represents significant events that impacts time series observations, and the error term t represents noisy data.

During model training, a time series can be divided into training and testing. The terms $g(t)$, $s(t)$ and $h(t)$ can be automatically inferred by classical statistical methods in the area of time series analysis, such as the Generalized Additive Model (GAM) (Hastie and Tibshirani, 1987) used in Prophet. In the training step, the terms are adjusted to find an additive model that best fits the known observations in the training time series. Next, we evaluated the model in new data, i.e., the testing time series.

In the case of the temporal dynamics of the software requirements, domain knowledge is represented by specific points (e.g. changepoints) in the time series that indicate potential growth of the requirement in negative reviews. Figure 6 shows the forecasting for a software requirement. Original observations are the black dots and the blue line represents the forecast model. The light blue area is the confidence interval of the predictions. The vertical dashed lines are the time series changepoints.

Changepoints play an important role in forecasting models, as they represent abrupt changes in the trend. Changepoints can be estimated automatically during model training, but domain knowledge, such as the date of app releases, marketing campaigns, and server failures, are changepoints that can be added

manually by software engineers. Therefore, the changepoints could be specified by the analyst using known dates of product launches and other growth-altering events or may be automatically selected given a set of candidates. In MAPP-Reviews, we have two possible options for selecting changepoints in the predictive model. The first option is automatic changepoint selection, where the Prophet specifies 25 potential changepoints which are uniformly placed in the first 80% of the time series. The second option is the manual specification which has a set of dates provided by a domain analyst. In this case, the changepoints could be entirely limited to a small set of dates. If no known dates are provided, by default we use the most recent observations which have a value greater than the average of the observations, i.e., we want to emphasize the highest peaks of the time series, as they indicate critical periods of negative revisions from the past.

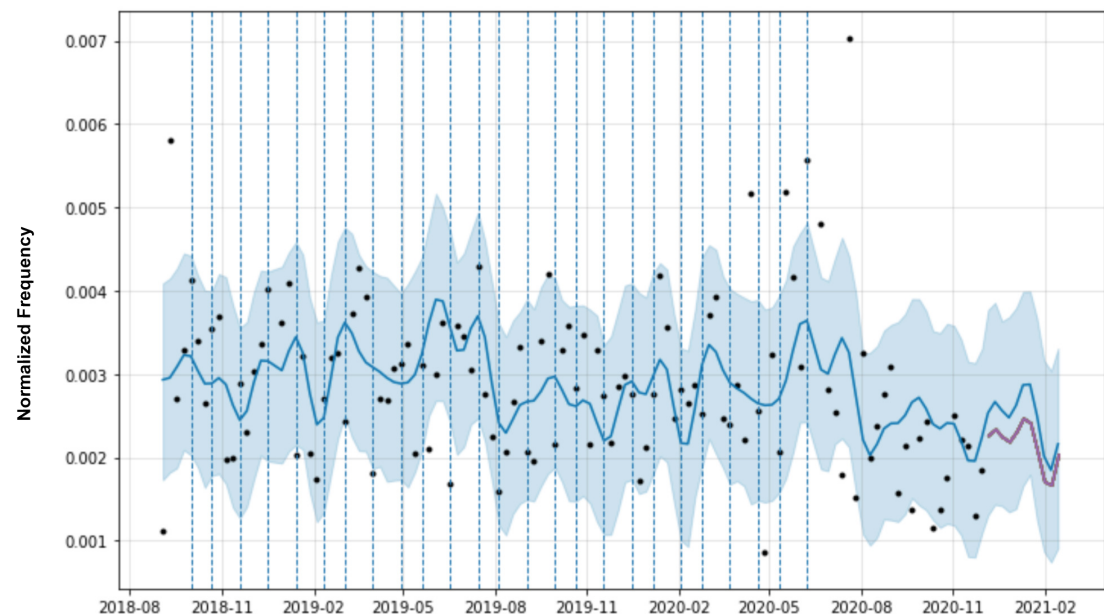


Figure 6. Prophet forecasting with automatic changepoints of a requirement.

In the experimental evaluation, we show the MAPP-Reviews ability to predict perceptually important points in the software requirements time series, allowing the identification of initial trends in defective requirements to support preventive strategies in software maintenance.

Table 4 shows an emerging issue being predicted 6 weeks in advance in the period from October 2020 to January 2021. The table presents a timeline represented by the horizon (h) in weeks, with the volume of negative raw reviews ($Vol.$). An example of a negative review is shown for each week until reaching the critical week (peak), with $h = 16$. The table row with $h = 10$ highlighted in bold shows when MAPP-Reviews identified the uptrend. In this case, we show the MAPP-Reviews alert for the “Time of arrival” requirement of the Uber Eats app. In particular, the emerging issue identified in the negative reviews is the low accuracy of the estimated delivery time in the app. The text of the user review samples has been entered in its entirety, without any pre-treatment. A graphical representation of this prediction is shown in Figure 7.

RESULTS

The proposed approach is validated through an experimental evaluation with popular food delivery apps. These apps represent a dynamic and complex environment consisting of restaurants, food consumers, and drivers operating in highly competitive conditions (Williams et al., 2020). In addition, this environment means a real scenario of commercial limitations, technological restrictions, and different user experience contexts, which makes detecting emerging issues early an essential task. For this experimental evaluation, we used a dataset with 86,610 reviews of three food delivery apps: Uber Eats, Foodpanda, and Zomato. The dataset was obtained in the first stage (App Reviews) of MAPP-Reviews and is available at

Table 4. Example of emerging issue prediction alert for the “Time of arrival” requirement of the Uber Eats app reviews triggered by MAPP-Reviews.

h	Vol.	Token	Review
1	768	Delivery time	Listed delivery times are inaccurate majority of the time.
2	849	Time frame	This app consistently gives incorrect, shorter delivery time frame to get you to order, but the deliveries are always late. The algorithm to predict the delivery time should be fixed so that you'll stop lying to your customers.
3	896	Arrival time	Ordered food and they told me it was coming. The wait time was supposed to be 45 minutes. They kept pushing back the arrival time, and we waited an hour and 45 minutes for food, only to have them CANCEL the order and tell us it wasn't coming. If an order is unable to be placed you need to tell customers BEFORE they've waited almost 2 HOURS for their food.
4	1247	Delivery time	The app was easy to navigate but the estimated delivery time kept changing and it took almost 2hrs to receive food and I live less than 4 blocks away pure ridiculousness if I would of know that I would of just walked there and got it.
5	1056	Estimated time	Everyone cancels and it ends up taking twice the estimated time to get the food delivered. You dont get updated on delays unless you actively monitor. Uber has failed at food delivery.
6	997	More time	Uber Eats lies. Several occasions showed delays because "the restaurant requested more time" but really it was Uber Eats unable to find a driver. I called the restaurants and they said the food has been ready for over an hour!
7	939	Delivery time	Your app is unintuitive. Delivery times are wildly inaccurate and orders are canceled with no explanation, information or help.
8	854	Estimated time	This service is terrible. Delivery people never arrive during the estimated time.
9	994	Time	Delivery times increase significantly once your order is accepted. 25-45 mins went up to almost 2hours! Not easy to cancel. Also one restaurant that looked available said I was too far away after I had filled my basket. Other than that the app is easy to use.
10	1257	Time estimate	Use door dash or post mates, uber eats has definitely gone down in quality. Extremely inaccurate time estimates and they ignore your support requests until its to late to cancel an order and get a refund.
11	1443	Delivery time	Delivery times are constantly updated, what was estimated at 25-35 minutes takes more than two hours. I understand it's just an estimate, but 4X that is ridiculous.
12	1478	Delivery time	Inaccurate delivery time
13	1376	Estimated time	Used to use this app a lot. Ever since they made it so you have to pay for your delivery to come on time the app is useless. You will be stuck waiting for food for an hour most of the time. The estimated time of arrival is never accurate. Have had my food brought to wrong addresses or not brought at all. I will just take the extra time out of my day to pick up the food myself rather than use this app.
14	1446	Estimated time	Terrible, the estimated time of arrival is never accurate and has regularly been up to 45 MINUTES LATE with no refund. Doordash is infinitely better, install that instead, it also has more restaurants
15	1354	Estimated time	App is good but this needs to be more reliable on its service. the estimated arrival time needs to be matched or there should be a option to cancel the order if they couldnt deliver on estimated time. Continuesly changing the estimated delivery time after the initial order confirmation is inappropriate.
16	1627	Estimated time	I use this app a lot and recently my order are always late at least double the time im originally quoted. Every time my food is cold. Maybe the estimated time should be adjusted to reflect what the actual time may be.

391 <https://github.com/vitormesaque/mapp-reviews>. The choice of these apps was based on their popularity
392 and the number of reviews available. The reviews are from September 2018 to January 2021.

393 After the software requirements extraction and clustering stage (with $k = 300$ clusters), the six most
394 popular (frequent) requirements clusters were considered for time series prediction. The following
395 software requirements clusters were selected: “Ordering”, “Go pick up”, “Delivery”, “Arriving time”,
396 “Advertising”, and “Payment”. The requirements clusters are shown in Table 5 with the associated words
397 ordered by silhouette.

398 In the MAPP-Reviews prediction stage, we evaluated two scenarios using Prophet. The first scenario
399 is the baseline, where we use the automatic parameters fitting of the Prophet. By default, Prophet will
400 automatically detect the changepoints. In the second scenario, we specify the potential changepoints,
401 thereby providing domain knowledge for software requirements rather than automatic changepoint
402 detection. Therefore, the changepoint parameters are used when we provide the dates of the changepoints
403 instead of the Prophet determining them. In this case, we use the most recent observations that have a
404 value greater than the average of observations, i.e., critical periods with high frequencies of negative
405 reviews in the past.

406 We used the MAPE (Mean Absolute Percentage Error) metric to evaluate the forecasting performance
407 (Makridakis, 1993), as defined in Equation 8,

Table 5. Software requirements clusters for food delivery apps used in the experimental evaluation. Tokens well allocated in each cluster (silhouette measure) were selected to support the cluster labeling.

Cluster Label	Tokens with Silhouette values (<i>s</i>)
Ordering	“ordering” (<i>s</i> = 0.1337), “order’s” (<i>s</i> = 0.1250), “order from” (<i>s</i> = 0.1243), “order will” (<i>s</i> = 0.1221), “order” (<i>s</i> = 0.1116), “the order”, (<i>s</i> = 0.1111)
Go pick up	“go pick up”(s = 0.1382), “pick up the” (<i>s</i> = 0.1289), “pick up at”, (<i>s</i> = 0.1261), “to take” (<i>s</i> = 0.1176), “go get” (<i>s</i> = 0.1159)
Delivery	“delivering parcels” (<i>s</i> = 0.1705), “delivery options” (<i>s</i> = 0.1590), “waive delivery” (<i>s</i> = 0.1566), “delivery charges” (<i>s</i> = 0.1501), “accept delivery” (<i>s</i> = 0.1492)
Arriving time	“arrival time” (<i>s</i> =0.3303), “waisting time” (<i>s</i> = 0.3046), “arriving time” (<i>s</i> = 0.3042), “estimate time” (<i>s</i> = 0.2877), “delievery time” (<i>s</i> = 0.2743)
Advertising	“anoyning ads” (<i>s</i> = 0.3464), “pop-up ads” (<i>s</i> = 0.3440), “ads pop up” (<i>s</i> = 0.3388), “commercials advertise” (<i>s</i> = 0.3272), “advertising” (<i>s</i> = 0.3241)
Payment	“payment getting” (<i>s</i> = 0.2618), “payment get” (<i>s</i> = 0.2547), “getting payment” (<i>s</i> = 0.2530), “take payment”(s = 0.2504), “payment taking” (<i>s</i> = 0.2471), “payment” (<i>s</i> = 0.2401)

$$MAPE = \frac{1}{h} \sum_{t=1}^h \frac{|real_t - pred_t|}{real_t} \quad (8)$$

where $real_t$ is the real value and $pred_t$ is the predicted value by the method, and h is the number of forecast observations in the estimation period (prediction horizon). In practical terms, MAPE is a measure of the percentage error that, in a simulation, indicates how close the prediction was made to the known values of the time series. We consider a prediction horizon (h) ranging from 1 to 4, with weekly seasonality.

Table 6 summarizes the main experimental results. The first scenario (1) with the default parameters obtains superior results compared to the second scenario (2) for all forecast horizons. In general, automatic changepoints obtains 9.33% of model improvement, considering the average of MAPE values from all horizons ($h = 1$ to $h = 4$).

Table 6. Comparison of MAPE in General.

h	MAPE (Mean ± SD)	
	(1) Automatic changepoint	(2) Specifying the changepoints
1	13.82 ± 16.42	15.47 ± 14.42
2	15.58 ± 19.09	16.94 ± 17.20
3	16.26 ± 20.18	17.60 ± 18.71
4	16.09 ± 19.24	17.47 ± 18.37

In particular, we are interested in the peaks of the series since our hypothesis is that the peaks represent potential problems in a given software requirement. Thus, Table 7 shows MAPE calculated only for time series peaks during forecasting. In this case, predictions with the custom changepotins locations (scenario 2) obtained better results than the automatic detection for all prediction horizons ($h = 1$ to $h = 4$), obtaining 3.82% of forecasting improvement. These results provide evidence that domain knowledge can improve the detection of potential software requirements to be analyzed for preventive maintenance.

In particular, analyzing the prediction horizon, the results show that the best predictions were obtained with $h = 1$ (1 week). In practical terms, this means the initial trend of a defective requirement can be identified one week in advance.

Finally, to exemplify MAPP-Reviews forecasting, Figure 7 shows the training data (Arriving time software requirement) represented as black dots and the forecast as a blue line, with upper and lower

Table 7. MAPE analysis (at the peaks of the time series) of each scenario considering the software requirements.

h	MAPE (Mean \pm SD)	
	(1) Automatic changepoint	(2) Specifying the changepoints
1	10.65 \pm 8.41	10.30 \pm 8.06
2	11.61 \pm 8.80	11.00 \pm 8.71
3	11.81 \pm 8.86	11.42 \pm 8.52
4	11.49 \pm 8.71	11.19 \pm 8.34

427 bounds in a blue shaded area. At the end of the time series, the darkest line is the real values plotted over
 428 the predicted values in blue. The lines plotted vertically represent the changepoints.

429 For reproducibility purposes, we provide a GitHub repository at [https://github.com/vitormesaque/mapp-](https://github.com/vitormesaque/mapp-reviews)
 430 reviews containing the source code and details of each stage of the method, as well as the raw data and all
 431 the results obtained.

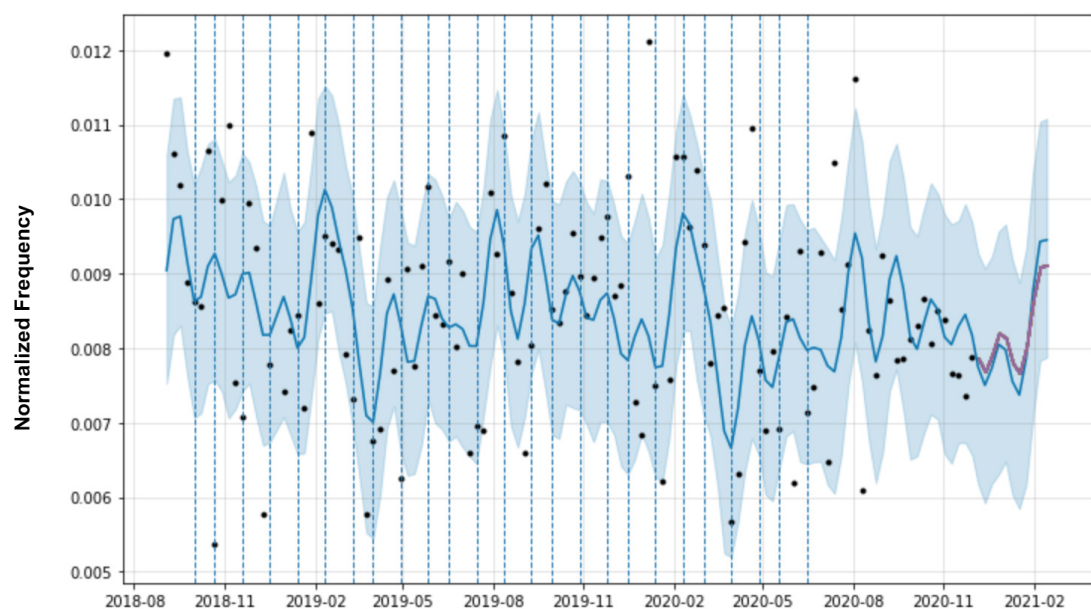


Figure 7. Forecasting for software requirement cluster (Arriving time) from Uber Eats App reviews.

432 DISCUSSION

433 Timely and effective detection of software requirements issues is crucial for app developers. The results
 434 show that MAPP-Reviews can detect significant points in the time series that provide information about
 435 the future behavior of a software requirement, allowing software engineers to anticipate the identification
 436 of emerging issues that may affect app evaluation. An issue related to a software requirement reported
 437 in user reviews is defined as an emerging issue when there is an upward trend for that requirement in
 438 negative reviews. Our method trains predictive models to identify requirements with higher negative
 439 evaluation trends, but a negative review will inevitably impact the rating. However, our objective is to
 440 mitigate this negative impact.

441 The prediction horizon (h) is an essential factor in detecting emerging issues to mitigate negative
 442 impacts. Software engineers and the entire development team need to know as soon as possible about
 443 software problems to anticipate them. In this context would not be feasible to predict the following months
 444 as it is tough to find a correlation between what happens today and what will happen in the next few

months about bug reports. Therefore, MAPP-Reviews forecasts at the week level. This strategy allows us to identify the issues that are starting to happen and predict whether they will worsen in the coming weeks. Even at the week level, the best forecast should be with the shortest forecast horizon, i.e., one week ($h = 1$). A longer horizon, i.e., three ($h = 3$) or four weeks ($h = 4$), could be too late to prevent an issue from becoming severe and having more impact on the overall app rating. The experimental evaluation shows that our method obtains the best predictions with the shortest horizon ($h = 1$). In practical terms, this means that MAPP-Reviews identifies the initial trend of a defective requirement a week in advance. In addition, we can note that a prediction error rate (MAPE) of up to 20% is acceptable. For example, consider that the prediction is 1000 negative reviews for a specific requirement at a given point, but the model predicts 800 negative reviews. Even with 20% of MAPE, we can identify a significant increase in negative reviews for a requirement and trigger alerts for preventive software maintenance, i.e., when MAPP-Reviews predicts an uptrend, the software development team should receive an alert. In the time series forecast shown in Figure 7, we observe that the model would be able to predict the peaks of negative reviews for the software requirement one week in advance.

The forecast presented in Figure 7 shows that the model was able to predict the peak of negative reviews for the “Arriving time” requirement. An emerging issue detection system based only on the frequency of a topic could trigger many false detections, i.e., it would not detect defective functionality but issues related to the quality of services offered. Analyzing user reviews, we found that some complaints are about service issues rather than defective requirements. For example, the user may complain about the delay in the delivery service and negatively rate the app, but in reality, they are complaining about the restaurant, i.e., a problem with the establishment service. We’ve seen that this pattern of user complaints is repeated across other app domains, not just the food delivery service. In delivery food apps, these complaints about service are constant, uniform, and distributed among all restaurants available in the app. In Table 4, it is clear that the emerging issue refers to the deficient implementation of the estimated delivery time prediction functionality. Our results show that when there is a problem in the app related to a defective software requirement, there are increasing complaints associated with negative reviews regarding that requirement.

An essential feature in MAPP-Reviews is changepoints. Assume that a time series represents the evolution of a software requirement over time, observing negative reviews for this requirement. Also, consider that time series frequently have abrupt changes in their trajectories. Given this, the changepoints describe abrupt changes in the time series trend, i.e., means a specific date that indicates a trend change. Therefore, specifying custom changepoints becomes significantly important for the predictive model because the uptrend in time series can also be associated with domain knowledge factors. By default, our model will automatically detect these changepoints. However, we have found that specifying custom changepoints improves prediction significantly in critical situations for the emerging issue detection problem. In general, the automatic detection of changepoints had better MAPE results in most evaluations. However, the custom changepoints obtained the best predictions at the time series peaks for all horizons ($h = 1$ to $h = 4$) of experiment simulations. Our experiment suggests a greater interest in identifying potential defective requirements trends in the time series peaks. As a result, we conclude that specifying custom changepoints in the predictive model is the best strategy to identify potential emerging issues.

Furthermore, the results indicate the potential impact of incorporating changepoints into the predictive model using the information of app developers, i.e., defining specific points over time with a meaningful influence on app evaluation. In addition, software engineers can provide sensitive company data and domain knowledge to explore and improve the predictive model potentially. For this purpose, we depend on sensitive company data related to the software development and management process, e.g., release planning, server failures, and marketing campaigns. In particular, we can investigate the relationship between the release dates of app updates and the textual content of the update publication with the upward trend in negative evaluations of a software requirement. In a real-world scenario in the industry, software engineers using MAPP-Reviews will provide domain-specific information.

We show that MAPP-Reviews provides software engineers with tools to perform software maintenance activities, particularly preventive maintenance, by automatically monitoring the temporal dynamics of software requirements.

The results of our research show there are new promising prospects for the future, and new possibilities for innovation research in this area emerge with our results so far. We intend to explore further our method to deeply determine the input variables that most contribute to the output behavior and the non-influential

inputs or to determine some interaction effects within the model. In addition, sensitivity analysis can help us reduce the uncertainties found more effectively and calibrate the model.

Limitations

Despite the significant results obtained, we can still improve the predictive model. In the scope of our experimental evaluation, we only investigate the incorporation of software domain-specific information through trend changepoints. Company-sensitive information and the development team's domain knowledge were not considered in the predictive model because we don't have access to this information. Therefore, we intend to evaluate our proposed method in the industry and explore more specifics of the domain knowledge to improve the predictive model.

Another issue that is important to highlight is the sentiment analysis in app reviews. We assume that it is possible to improve the classification of negative reviews by incorporating sentiment analysis techniques. We can incorporate a polarity classification stage (positive, negative, and neutral) of the extracted requirement, allowing a software requirements-based sentiment analysis. In the current state of our research, we only consider negative reviews with low ratings and associate them with the software requirements mentioned in the review.

Finally, to use MAPP-Reviews in a real scenario, there must be already a sufficient amount of reviews distributed over time, i.e., a minimum number of time-series observations available for the predictive model to work properly. Therefore, in practical terms, our method is more suitable when large volumes of app reviews are available to be analyzed.

CONCLUSIONS

Opinion mining for app reviews can provide useful user feedback to support software engineering activities. We introduced the temporal dynamics of requirements analysis to predict initial trends on defective requirements from users' opinions before negatively impacting the overall app's evaluation. We presented the MAPP-Reviews (Monitoring App Reviews) approach to (1) extract and cluster software requirements, (2) generate time series with the time dynamics of requirements, (3) identify requirements with higher trends of negative evaluation.

The experimental results show that our method is able to find significant points in the time series that provide information about the future behavior of a requirement through app reviews, thereby allowing software engineers to anticipate the identification of requirements that may affect the app's evaluation. In addition, we show that it's beneficial to incorporate changepoints into the predictive model by using domain knowledge, i.e., defining points over time with significant impacts on the app's evaluation.

We compared the MAPP-Reviews in two scenarios: first using automatic changepoint detection and second specifying the changepoint locations. In particular, the automatic detection of points of change had better MAPE results in most evaluations. On the other hand, the best predictions at the time series peaks (where there is a greater interest in identifying potential defective requirements trends) were obtained by specifying changepoints.

Future work directions involve evaluating MAPP-Reviews in other scenarios to incorporate and compare several other types of domain knowledge into the predictive model, such as new app releases, marketing campaigns, server failures, competing apps, among other information that may impact the evaluation of apps. Another direction for future work is to implement a dashboard tool for monitoring app reviews, thus allowing the dispatching of alerts and reports.

REFERENCES

- AlSubaih, Afnan, Federica Sarro, Sue Black, Licia Capra, and Mark Harman (2019). "App store effects on software engineering practices". In: *IEEE Transactions on Software Engineering*.
- April, Alain and Alain Abran (2012). *Software maintenance management: evaluation and continuous improvement*. Vol. 67. John Wiley & Sons.
- Araujo, Adailton, Marcos Golo, Breno Viana, Felipe Sanches, Roseli Romero, and Ricardo Marcacini (2020). "From Bag-of-Words to Pre-trained Neural Language Models: Improving Automatic Classification of App Reviews for Requirements Engineering". In: *Anais do XVII Encontro Nacional de Inteligência Artificial e Computacional*. SBC, pp. 378–389.

- 550 Araujo, Adailton and Ricardo Marcondes Marcacini (2021). “RE-BERT: Automatic Extraction of Soft-
551 ware Requirements from App Reviews using BERT Language Model”. In: *The 36th ACM/SIGAPP
552 Symposium On Applied Computing*. DOI: 10.1145/3412841.3442006.
- 553 Bennett, Keith H. and Václav T. Rajlich (2000). “Software Maintenance and Evolution: A Roadmap”. In:
554 *Proceedings of the Conference on The Future of Software Engineering*. ICSE ’00. Association for
555 Computing Machinery: Limerick, Ireland, pp. 73–87. DOI: 10.1145/336512.336534.
- 556 Blei, David M., Andrew Y. Ng, and Michael I. Jordan (Mar. 2003). “Latent Dirichlet Allocation”. In: *J.
557 Mach. Learn. Res.* 3(null), pp. 993–1022. ISSN: 1532-4435.
- 558 Bojanowski, Piotr, Edouard Grave, Armand Joulin, and Tomas Mikolov (2017). *Enriching Word Vectors
559 with Subword Information*. arXiv: 1607.04606 [cs.CL].
- 560 Chatfield, Chris and Haipeng Xing (2019). *The analysis of time series: an introduction with R*. CRC press.
- 561 Chen, Miao and Xiaozong Liu (2011). “Predicting Popularity of Online Distributed Applications: iTunes
562 App Store Case Analysis”. In: *Proceedings of the 2011 IConference*. iConference ’11. Association
563 for Computing Machinery: Seattle, Washington, USA, pp. 661–663. DOI: 10.1145/1940761.
564 1940859.
- 565 Chen, Ning, Jialiu Lin, Steven CH Hoi, Xiaokui Xiao, and Boshen Zhang (2014). “AR-miner: mining
566 informative reviews for developers from mobile app marketplace”. In: *Proceedings of the 36th
567 international conference on software engineering*, pp. 767–778.
- 568 Ciurumelea, Adelina, Andreas Schaufelbühl, Sebastiano Panichella, and Harald C. Gall (2017). “Analyz-
569 ing reviews and code of mobile apps for better release planning”. In: *2017 IEEE 24th International
570 Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 91–102.
- 571 Dabrowski, Jacek, Emmanuel Letier, Anna Perini, and Angelo Susi (2020). “Mining User Opinions
572 to Support Requirement Engineering: An Empirical Study”. In: *Advanced Information Systems
573 Engineering*. Ed. by Schahram Dustdar, Eric Yu, Camille Salinesi, Dominique Rieu, and Vik Pant.
574 Springer International Publishing: Cham, pp. 401–416.
- 575 Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). “BERT: Pre-training of
576 Deep Bidirectional Transformers for Language Understanding”. In: *CoRR* abs/1810.04805. arXiv:
577 1810.04805.
- 578 Di Sorbo, Andrea, Sebastiano Panichella, Carol V. Alexandru, Junji Shimagaki, Corrado A. Visaggio,
579 Gerardo Canfora, and Harald C. Gall (2016). “What Would Users Change in My App? Summarizing
580 App Reviews for Recommending Software Changes”. In: *Proceedings of the 2016 24th ACM SIG-
581 SOFT International Symposium on Foundations of Software Engineering*. FSE 2016. Association for
582 Computing Machinery: Seattle, WA, USA, pp. 499–510. ISBN: 9781450342186.
- 583 Dragoni, Mauro, Marco Federici, and Andi Rexha (2019). “An unsupervised aspect extraction strategy for
584 monitoring real-time reviews stream”. In: *Information Processing Management* 56(3), pp. 1103–1118.
585 DOI: <https://doi.org/10.1016/j.ipm.2018.04.010>.
- 586 Galvis Carreño, Laura V. and Kristina Winbladh (2013). “Analysis of User Comments: An Approach for
587 Software Requirements Evolution”. In: *Proceedings of the 2013 International Conference on Software
588 Engineering*. ICSE ’13. IEEE Press: San Francisco, CA, USA, pp. 582–591.
- 589 Gao, Cuiyun, Jichuan Zeng, Michael R. Lyu, and Irwin King (2018). “Online App Review Analysis
590 for Identifying Emerging Issues”. In: *Proceedings of the 40th International Conference on Software
591 Engineering*. ICSE ’18. Association for Computing Machinery: Gothenburg, Sweden, pp. 48–58. DOI:
592 10.1145/3180155.3180218.
- 593 Gao, Cuiyun, Jichuan Zeng, Zhiyuan Wen, David Lo, Xin Xia, Irwin King, and Michael R. Lyu (2020).
594 *Emerging App Issue Identification via Online Joint Sentiment-Topic Tracing*. arXiv: 2008.09976
595 [cs.SE].
- 596 Gao, Cuiyun, Wujie Zheng, Yuetang Deng, David Lo, Jichuan Zeng, Michael R. Lyu, and Irwin King
597 (2019). “Emerging App Issue Identification from User Feedback: Experience on WeChat”. In: *2019
598 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice
599 (ICSE-SEIP)*, pp. 279–288. DOI: 10.1109/ICSE-SEIP.2019.00040.
- 600 Gómez, Maria, Romain Rouvoy, Martin Monperrus, and Lionel Seinturier (2015). “A Recommender
601 System of Buggy App Checkers for App Store Moderators”. In: *Proceedings of the Second ACM
602 International Conference on Mobile Software Engineering and Systems*. MOBILESoft ’15. IEEE
603 Press: Florence, Italy, pp. 1–11.

- 604 Gu, Xiaodong and Sunghun Kim (2015). ““What Parts of Your Apps are Loved by Users?” (T)”. In: *2015*
605 *30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pp. 760–770.
- 606 Guzman, Emitza, Rana Alkadhi, and Norbert Seyff (2016). “A Needle in a Haystack: What Do Twitter
607 Users Say about Software?” In: *2016 IEEE 24th International Requirements Engineering Conference*
608 *(RE)*, pp. 96–105.
- 609 Guzman, Emitza, Rana Alkadhi, and Norbert Seyff (Sept. 2017). “An Exploratory Study of Twitter
610 Messages about Software Applications”. In: *Requir. Eng.* 22(3), pp. 387–412. ISSN: 0947-3602.
- 611 Guzman, Emitza and Walid Maalej (2014). “How do users like this feature? a fine grained sentiment
612 analysis of app reviews”. In: *2014 IEEE 22nd international requirements engineering conference*
613 *(RE)*. IEEE, pp. 153–162.
- 614 Harman, Mark, Yue Jia, and Yuanyuan Zhang (2012). “App Store Mining and Analysis: MSR for App
615 Stores”. In: *Proceedings of the 9th IEEE Working Conference on Mining Software Repositories*. MSR
616 ’12. IEEE Press: Zurich, Switzerland, pp. 108–111.
- 617 Hastie, Trevor and Robert Tibshirani (1987). “Generalized additive models: some applications”. In:
618 *Journal of the American Statistical Association* 82(398), pp. 371–386.
- 619 Hoon, Leonard, Rajesh Vasa, Jean-Guy Schneider, and Kon Mouzakis (2012). “A Preliminary Analysis of
620 Vocabulary in Mobile App User Reviews”. In: *Proceedings of the 24th Australian Computer-Human*
621 *Interaction Conference*. OzCHI ’12. Association for Computing Machinery: Melbourne, Australia,
622 pp. 245–248. DOI: 10.1145/2414536.2414578.
- 623 Iacob, Claudia and Rachel Harrison (2013). “Retrieving and Analyzing Mobile Apps Feature Requests
624 from Online Reviews”. In: *Proceedings of the 10th Working Conference on Mining Software Reposito-*
625 *ries*. MSR ’13. IEEE Press: San Francisco, CA, USA, pp. 41–44.
- 626 Johann, Timo, Christoph Stanik, Walid Maalej, et al. (2017). “Safe: A simple approach for feature
627 extraction from app descriptions and app reviews”. In: *2017 IEEE 25th International Requirements*
628 *Engineering Conference (RE)*. IEEE, pp. 21–30.
- 629 Johanssen, Jan Ole, Anja Kleebaum, Bernd Bruegge, and Barbara Paech (2019). “How do practitioners
630 capture and utilize user feedback during continuous software engineering?” In: *2019 IEEE 27th*
631 *International Requirements Engineering Conference (RE)*. IEEE, pp. 153–164.
- 632 Khalid, Hammad, Emad Shihab, Meiyappan Nagappan, and Ahmed E. Hassan (2015). “What Do Mobile
633 App Users Complain About?” In: *IEEE Software* 32(3), pp. 70–77. DOI: 10.1109/MS.2014.50.
- 634 Licorish, Sherlock A., Bastin Tony Roy Savarimuthu, and Swetha Keertipati (2017). “Attributes That
635 Predict Which Features to Fix: Lessons for App Store Mining”. In: *Proceedings of the 21st Interna-*
636 *tional Conference on Evaluation and Assessment in Software Engineering*. EASE’17. Association for
637 Computing Machinery: Karlskrona, Sweden, pp. 108–117. DOI: 10.1145/3084226.3084246.
- 638 Lientz, Bennett P and E Burton Swanson (1980). *Software maintenance management*. Addison-Wesley
639 Longman Publishing Co., Inc.
- 640 Liu, Bing (2012). “Sentiment analysis and opinion mining”. In: *Synthesis lectures on human language*
641 *technologies* 5(1), pp. 1–167.
- 642 Maalej, Walid and Hadeer Nabil (2015). “Bug report, feature request, or simply praise? On automatically
643 classifying app reviews”. In: *2015 IEEE 23rd International Requirements Engineering Conference*
644 *(RE)*, pp. 116–125.
- 645 MacQueen, James et al. (1967). “Some methods for classification and analysis of multivariate observa-
646 tions”. In: *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*.
647 Vol. 1. 14. Oakland, CA, USA, pp. 281–297.
- 648 Makridakis, Spyros (1993). “Accuracy measures: theoretical and practical concerns”. In: *International*
649 *Journal of Forecasting* 9(4), pp. 527–529. DOI: [https://doi.org/10.1016/0169-](https://doi.org/10.1016/0169-2070(93)90079-3)
650 [2070\(93\)90079-3](https://doi.org/10.1016/0169-2070(93)90079-3).
- 651 Malik, Haroon, Elhadi M. Shakshuki, and Wook-Sung Yoo (2020). “Comparing mobile apps by identifying
652 ‘Hot’ features”. In: *Future Generation Computer Systems* 107, pp. 659–669. ISSN: 0167-739X. DOI:
653 <https://doi.org/10.1016/j.future.2018.02.008>.
- 654 Martin, William, Federica Sarro, Yue Jia, Yuanyuan Zhang, and Mark Harman (2016). “A survey of
655 app store analysis for software engineering”. In: *IEEE transactions on software engineering* 43(9),
656 pp. 817–847.
- 657 Mcilroy, Stuart, Nasir Ali, and Ahmed E. Hassan (June 2016). “Fresh Apps: An Empirical Study of
658 Frequently-Updated Mobile Apps in the Google Play Store”. In: *Empirical Softw. Engg.* 21(3),

- pp. 1346–1370. ISSN: 1382-3256. DOI: 10.1007/s10664-015-9388-2. URL: <https://doi.org/10.1007/s10664-015-9388-2>.
- Mcilroy, Stuart, Nasir Ali, Hammad Khalid, and Ahmed E. Hassan (June 2016). “Analyzing and Automatically Labelling the Types of User Issues That Are Raised in Mobile App Reviews”. In: *Empirical Softw. Engg.* 21(3), pp. 1067–1106. DOI: 10.1007/s10664-015-9375-7.
- Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean (2013). *Efficient Estimation of Word Representations in Vector Space*. arXiv: 1301.3781 [cs.CL].
- Nayebi, Maleknaz, Bram Adams, and Guenther Ruhe (2016). “Release Practices for Mobile Apps – What do Users and Developers Think?” In: *2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER)*. Vol. 1, pp. 552–562. DOI: 10.1109/SANER.2016.116.
- Nayebi, Maleknaz, Henry Cho, and Guenther Ruhe (Oct. 2018). “App Store Mining is Not Enough for App Improvement”. In: *Empirical Softw. Engg.* 23(5), pp. 2764–2794. ISSN: 1382-3256.
- Nayebi, Maleknaz, Konstantin Kuznetsov, Paul Chen, Andreas Zeller, and Guenther Ruhe (2018). “Anatomy of Functionality Deletion: An Exploratory Study on Mobile Apps”. In: *Proceedings of the 15th International Conference on Mining Software Repositories*. MSR ’18. Association for Computing Machinery: Gothenburg, Sweden, pp. 243–253. ISBN: 9781450357166.
- Nayebi, Maleknaz, Mahshid Marbouti, Rache Quapp, Frank Maurer, and Guenther Ruhe (2017). “Crowd-sourced Exploration of Mobile App Features: A Case Study of the Fort McMurray Wildfire”. In: *2017 IEEE/ACM 39th International Conference on Software Engineering: Software Engineering in Society Track (ICSE-SEIS)*, pp. 57–66.
- Noei, Ehsan, Feng Zhang, and Ying Zou (2021). “Too Many User-Reviews! What Should App Developers Look at First?” In: *IEEE Transactions on Software Engineering* 47(2), pp. 367–378. DOI: 10.1109/TSE.2019.2893171.
- Pagano, D. and W. Maalej (2013). “User feedback in the appstore: An empirical study”. In: *2013 21st IEEE International Requirements Engineering Conference (RE)*, pp. 125–134. DOI: 10.1109/RE.2013.6636712.
- Palomba, Fabio, Mario Linares-Vásquez, Gabriele Bavota, Rocco Oliveto, Massimiliano Di Penta, Denys Poshyvanyk, and Andrea De Lucia (2015). “User reviews matter! Tracking crowdsourced reviews to support evolution of successful apps”. In: *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pp. 291–300.
- Palomba, Fabio, Mario Linares-Vásquez, Gabriele Bavota, Rocco Oliveto, Massimiliano Di Penta, Denys Poshyvanyk, and Andrea De Lucia (2018). “Crowdsourcing user reviews to support the evolution of mobile apps”. In: *Journal of Systems and Software* 137, pp. 143–162. ISSN: 0164-1212. DOI: <https://doi.org/10.1016/j.jss.2017.11.043>. URL: <https://www.sciencedirect.com/science/article/pii/S0164121217302807>.
- Panichella, Sebastiano, Andrea Di Sorbo, Emitza Guzman, Corrado A. Visaggio, Gerardo Canfora, and Harald C. Gall (2015). “How can i improve my app? Classifying user reviews for software maintenance and evolution”. In: *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pp. 281–290.
- Panichella, Sebastiano, Andrea Di Sorbo, Emitza Guzman, Corrado A. Visaggio, Gerardo Canfora, and Harald C. Gall (2016). “ARdoc: App Reviews Development Oriented Classifier”. In: *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*. FSE 2016. Association for Computing Machinery: Seattle, WA, USA, pp. 1023–1027. ISBN: 9781450342186.
- Pennington, Jeffrey, Richard Socher, and Christopher Manning (Oct. 2014). “GloVe: Global Vectors for Word Representation”. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics: Doha, Qatar, pp. 1532–1543. DOI: 10.3115/v1/D14-1162.
- Phong, Minh Vu, Tam The Nguyen, Hung Viet Pham, and Tung Thanh Nguyen (2015). “Mining User Opinions in Mobile App Reviews: A Keyword-Based Approach (T)”. In: *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pp. 749–759. DOI: 10.1109/ASE.2015.85.
- Reimers, Nils and Iryna Gurevych (2019). “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*

- 713 *Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-*
714 *IJCNLP)*, pp. 3973–3983.
- 715 Rousseeuw, Peter J (1987). “Silhouettes: a graphical aid to the interpretation and validation of cluster
716 analysis”. In: *Journal of computational and applied mathematics* 20, pp. 53–65.
- 717 Saaty, T.L. (1980). *The Analytic Hierarchy Process: Planning, Priority Setting, Resource Allocation*.
718 Advanced book program. McGraw-Hill International Book Company.
- 719 Sarro, Federica, Afnan A. Al-Subaih, Mark Harman, Yue Jia, William Martin, and Yuanyuan Zhang
720 (2015). “Feature lifecycles as they spread, migrate, remain, and die in App Stores”. In: *2015 IEEE*
721 *23rd International Requirements Engineering Conference (RE)*, pp. 76–85.
- 722 Al-Subaih, Afnan A, Federica Sarro, Sue Black, Licia Capra, Mark Harman, Yue Jia, and Yuanyuan
723 Zhang (2016). “Clustering mobile apps based on mined textual features”. In: *Proceedings of the 10th*
724 *ACM/IEEE international symposium on empirical software engineering and measurement*, pp. 1–10.
- 725 Taylor, Sean J. and Benjamin Letham (Jan. 2018). “Forecasting at Scale”. In: *The American Statistician*
726 72(1), pp. 37–45.
- 727 Tudor, D. and G.A. Walter (2006). “Using an agile approach in a large, traditional organization”. In:
728 *AGILE 2006 (AGILE’06)*, 7 pp.–373. DOI: 10.1109/AGILE.2006.60.
- 729 Vasa, Rajesh, Leonard Hoon, Kon Mouzakis, and Akihiro Noguchi (2012). “A Preliminary Analysis of
730 Mobile App User Reviews”. In: *Proceedings of the 24th Australian Computer-Human Interaction*
731 *Conference. OzCHI ’12*. Association for Computing Machinery: Melbourne, Australia, pp. 241–244.
732 DOI: 10.1145/2414536.2414577.
- 733 Vendramin, Lucas, Ricardo JGB Campello, and Eduardo R Hruschka (2010). “Relative clustering validity
734 criteria: A comparative overview”. In: *Statistical analysis and data mining: the ASA data science*
735 *journal* 3(4), pp. 209–235.
- 736 Villaruel, Lorenzo, Gabriele Bavota, Barbara Russo, Rocco Oliveto, and Massimiliano Di Penta (2016).
737 “Release Planning of Mobile Apps Based on User Reviews”. In: *2016 IEEE/ACM 38th International*
738 *Conference on Software Engineering (ICSE)*, pp. 14–24.
- 739 Vu, Phong Minh, Hung Viet Pham, Tam The Nguyen, and Tung Thanh Nguyen (2016). “Phrase-based
740 extraction of user opinions in mobile app reviews”. In: *2016 31st IEEE/ACM International Conference*
741 *on Automated Software Engineering (ASE)*, pp. 726–731.
- 742 Williams, Grant, Miroslav Tushev, Fahimeh Ebrahimi, and Anas Mahmoud (Dec. 2020). “Modeling user
743 concerns in Sharing Economy: the case of food delivery apps”. In: *Automated Software Engineering*
744 27. DOI: 10.1007/s10515-020-00274-7.
- 745 Zhang, Ying and Zhijie Lin (2018). “Predicting the helpfulness of online product reviews: A multilingual
746 approach”. In: *Electronic Commerce Research and Applications* 27, pp. 1–10.
- 747 Zhao, Liping, Waad Alhoshan, Alessio Ferrari, Keletso J. Letsholo, Muideen A. Ajagbe, Erol-Valeriu
748 Chioasca, and Riza T. Batista-Navarro (2020). *Natural Language Processing (NLP) for Requirements*
749 *Engineering: A Systematic Mapping Study*. arXiv: 2004.01099 [cs.SE].