ENSF 612: Fall 2021 Lecture - Course Project Overview

Dr. Gias Uddin, Assistant Professor Electrical and Software Engineering Schulich School of Engineering University of Calgary

https://giasuddin.ca/

Develop Supervised ML Models to Automatically Classify and Summarize Big Software Engineering Data

OPINER

A summarization engine for APIs

Search API

Explore review summaries

Search API Aspect

Find top ranked APIs by aspect

Search API Usage

Find Usage Scenarios For API

Most reviewed APIs

- 1. com.fasterxml.jackson
- 2. com.google.code.gson
- 3. org.springframework
- 4. org.glassfish.jersey
- 5. org.json
- 6. javax.xml
- 7. net.sf.json-lib
- 8. commons-httpclient
- 9. org.mongodb
- 10. com.google.gwt

Most reviewed aspects

- Usability
- 2. features
- 3. Bug
- 4. OnlySentiment
- 5. Documentation
- 6. Performance
- 7. Community
- 8. class
- 9. json
- 10. Security

http://opiner.polymtl.ca/

Most Used APIs

- 1. org.json
- 2. com.google.code.gson
- 3. com.fasterxml.jackson
- 4. javax.ws
- 5. org.springframework
- 6. commons-httpclient
- 7. net.sf.json-lib
- 8. org.glassfish.jersey
- 9. genson
- 10. xstream

Sources of big data - Software

Stack Exchange Online Technical Forums

Across all of Stack Overflow and the Stack Exchange network, we saw 9+ billion pageviews from 100+ million users over the course of the year.



JANUARY 18, 2019

State of the Stack 2019: A Year in Review

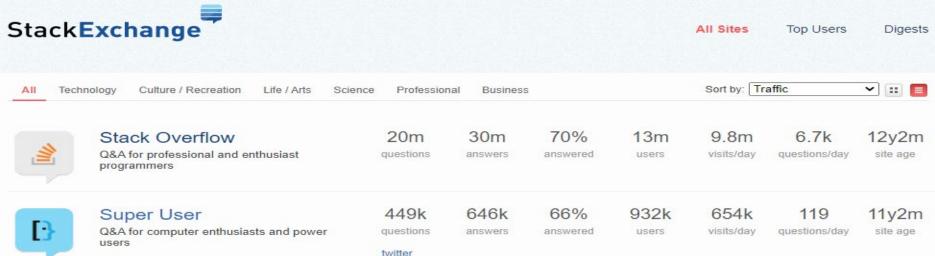
A loooong time ago, we used to post an annual "State of the Stack" update on the company and community. Then at some point it became an infographic which was... listen, everyone was doing infographics in 2011. Now it's 2019, the company has grown and changed in so many ways, so we're bringing this tradition...



David Fullerton

President and Chief Technology Officer (former)





Sources of big data - Software 100+ millions GitHub software repositories

40_{rm+}

developers on GitHub, including 10M new users in 2019.*

87 m+

pull requests merged in the last year—and 28% more developers opened their first pull request in 2019 than in 2018.*

44 m+

repositories created in the last year—and 44% more developers created their first repository in 2019 than in 2018.*

20 m+

issues closed in the last year. That's a lot of decisions made, bugs fixed, and boxes checked.*

Project Steps Per Group

- 1. Pick a research paper that labeled/summarized software engineering dataset
- Collect their dataset that they shared and that they used for labeling the dataset
- 3. Add 1000 more records into the dataset (details later)
- 4. Replicate the developed ML models in the paper using the original + extended dataset.
- 5. Extend the papers by experimenting with additional ML models
- 6. Show how the findings from ML models can be used to summarize the datasets and similar data in software engineering, e.g., by following Opiner

Categorizing the Content of GitHub README Files

Gede Artha Azriadi Prana · Christoph Treude · Ferdian Thung · Thushari Atapattu · David Lo

Received: date / Accepted: date

Abstract README files play an essential role in shaping a developer's first impression of a software repository and in documenting the software project that the repository hosts. Yet, we lack a systematic understanding of the content of a typical README file as well as tools that can process these files automatically. To close this gap, we conduct a qualitative study involving the manual annotation of 4,226 README file sections from 393 randomly sampled GitHub repositories and we design and evaluate a classifier and a set of features that can categorize these sections automatically. We find that information discussing the 'What' and 'How' of a repository is very common,

https://arxiv.org/pdf/1802.06997.pdf

Automating Intention Mining

Qiao Huang, Xin Xia, David Lo, Gail C. Murphy

Abstract—Developers frequently discuss aspects of the systems they are developing online. The comments they post to discussions form a rich information source about the system. Intention mining, a process introduced by Di Sorbo et al., classifies sentences in developer discussions to enable further analysis. As one example of use, intention mining has been used to help build various recommenders for software developers. The technique introduced by Di Sorbo et al. to categorize sentences is based on linguistic patterns derived from two projects. The limited number of data sources used in this earlier work introduces questions about the comprehensiveness of intention categories and whether the linguistic patterns used to identify the categories are generalizable to developer discussion recorded in other kinds of software artifacts (e.g., issue reports).

To assess the comprehensiveness of the previously identified intention categories and the generalizability of the linguistic patterns for category identification, we manually created a new dataset, categorizing 5,408 sentences from issue reports of four projects in GitHub. Based on this manual effort, we refined the previous categories. We assess Di Sorbo et al.'s patterns on this dataset, finding that the accuracy rate achieved is low (0.31). To address the deficiencies of Di Sorbo et al.'s patterns, we propose and investigate a convolution neural network (CNN)-based approach to automatically classify sentences into different categories of intentions. Our approach optimizes CNN by integrating batch normalization to accelerate the training speed, and an automatic hyperparameter tuning approach to tune appropriate hyperparameters of CNN. Our approach achieves an accuracy of 0.84 on the new dataset, improving Di Sorbo et al.'s approach by 171%. We also apply our approach to improve an automated software engineering task, in which we use our proposed approach to rectify misclassified issue reports, thus reducing the bias introduced by such data to other studies. A case study on four open source projects with 2,076 issue reports shows that our approach achieves an average AUC score of 0.687, which improves other baselines by at least 16%.

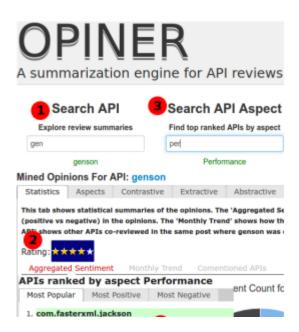
https://xin-xia.github.io/publication/tse185.pdf

Automatic Summarization of API Reviews

Gias Uddin School of Computer Science McGill University, Montréal, QC, Canada gias@cs.mcgill.ca Foutse Khomh SWAT Lab Polytechnique Montréal, QC, Canada foutse.khomh@polymtl.ca

Abstract—With the proliferation of online developer forums as informal documentation, developers often share their opinions about the APIs they use. However, given the plethora of opinions available for an API in various online developer forums, it can be challenging for a developer to make informed decisions about the APIs. While automatic summarization of opinions have been explored for other domains (e.g., cameras, cars), we found little research that investigates the benefits of summaries of public API reviews. In this paper, we present two algorithms (statistical and aspect-based) to summarize opinions about APIs. To investigate the usefulness of the techniques, we developed, Opiner, an online opinion summarization engine that presents summaries of opinions using both our proposed techniques and existing six off-theshelf techniques. We investigated the usefulness of Opiner using two case studies, both involving professional software engineers. We found that developers were interested to use our proposed summaries much more frequently than other summaries (daily vs once a year) and that while combined with Stack Overflow, Opiner helped developers to make the right decision with more accuracy and confidence and in less time.

Index Terms—Opinion mining, API informal documentation, opinion summaries, study, summary quality.



https://giasuddin.files.wordpress.com/2020/12/ase2017a-automatic-summarization-of-api-reviews.pdf

Sentiment Polarity Detection for Software Development

Fabio Calefato¹ · Filippo Lanubile² · Federico Maiorano² · Nicole Novielli²

Published online: 19 September 2017

© Springer Science+Business Media, LLC 2017

Abstract The role of sentiment analysis is increasingly emerging to study software developers' emotions by mining crowd-generated content within social software engineering tools. However, off-the-shelf sentiment analysis tools have been trained on non-technical domains and general-purpose social media, thus resulting in misclassifications of technical jargon and problem reports. Here, we present Senti4SD, a classifier specifically trained to support sentiment analysis in developers' communication channels. Senti4SD is trained and validated using a gold standard of Stack Overflow questions, answers, and comments manually annotated for sentiment polarity. It exploits a suite of both lexicon- and keyword-based features, as well as semantic features based on word embedding. With respect to a mainstream off-the-shelf tool, which we use as a baseline, Senti4SD reduces the misclassifications of neutral and positive posts as emotionally negative. To encourage replications, we release a lab package including the classifier, the word embedding space, and the gold standard with annotation guidelines.

https://arxiv.org/ftp/arxiv/papers/1709/1709.02984.pdf

Project Steps Per Group – Idea

- 1. Additional project ideas will be appreciated. Teams with new project ideas will be given higher preference
- 2. If more than one team picks same project idea, the teams need to differ from each other on the following key aspects:
 - The additional data to label (each team should label data different from each other)
 - The ML models that will be experimented with
- 3. If two teams are using the same project idea, the team with higher data quality and higher ML model performance will be given more grade than the other team.

Project Steps Per Group – Data Labeling Instructions

- 1. Each team will add at least 1000 new records to an existing dataset.
- 2. Each team will label the new 1000 records following approaches similar to the paper as described in the paper
- 3. Each team member will label the 1000 records separately as follows:
 - 1. Each member will label each of 1000 records
 - 2. For each record, the final label would be the one that most agree (e.g., two say I1 and one say I2, then the label is I1)
 - 3. For each record, if no such majority is found, a final label will be picked based on further discussions among the members