

Bug Report Prioritization

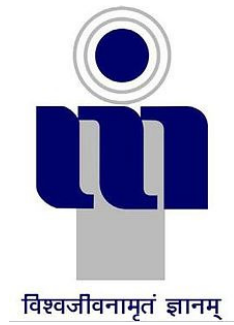
By

Anurag Yadav : 2020BCS-013

Under the Supervision of

Dr. Santosh Singh Rathore

Department of Computer Science



ABV-INDIAN INSTITUTE OF INFORMATION TECHNOLOGY
AND MANAGEMENT GWALIOR
GWALIOR, INDIA

DECLARATION

I hereby certify that the work, which is being presented in the report, entitled Bug Report Prioritization, in fulfillment of the requirement for the Colloquium Based on Summer Internship (BCCS-4105) and submitted to the institution is an authentic record of my/our own work carried out during the period May-2023 to July-2023 under the supervision of Dr. Santosh Singh Rathore. I also cited the reference about the text(s)/figure(s)/table(s) from where they have been taken.

Dated:

Signature of the candidate

This is to certify that the above statement made by the candidates is correct to the best of my knowledge.

Dated:

Signature of supervisor

Acknowledgements

I would like to express my sincere thanks to all those people who made this seminar report possible. First and foremost, I would like to express my profound respect and gratitude to my supervisor, Dr. Santosh Singh Rathore, who has been guiding force behind this work. I am greatly indebted for his invaluable guidance, constant encouragement, and his valuable comments on my work. I am fortunate enough to have such an advisor who gave me the freedom to think independently and explore new ideas. More importantly, I would like to thank for the patience he has shown in carefully reading and commenting on the manuscripts, and countless revisions of this dissertation. His commitments and dedication to research have been and will continue to be a constant source of inspiration for me. I have no doubts that finishing my degree in proper and timely manner was impossible without his help. I am highly privilege to have got an opportunity to work with such a wonderful person.

I am thankful to ABV IIITM for providing the research scholarship to undertake my research work. Finally, I would like to thank my friends and parents for the constant support towards me.

Anurag Yadav

Abstract

Software consumption has increased a lot in recent years. These software systems receive lots of bug reports. A Traigger is a person or program who reads he report and classifies the bugs and assigns priorities to these bugs so that bugs which are necessary to be fixed urgently could be fixed on time. Manually assigning priority to bugs is a time consuming task. Previously several other Machine Learning and Deep Learning techniques have been employed for the task of Bug Report Prioritization. Some models gave poor results whereas some models gave better results but their result were found to be inconsistent on newer Bug Reports. Thus, this project focuses on the prioritization of software bug reports using a Hierarchical Attention Network (HAN) model. The objective is to effectively handle class imbalance in bug report data and improve classification accuracy. The bug reports are preprocessed through text normalization, tokenization using the DistilBERT tokenizer, and attention mechanisms to capture hierarchical relationships within the text. The model's architecture incorporates bidirectional GRU layers and attention mechanisms at both word and sentence levels. The implementation includes data preprocessing, resampling using ADASYN for class imbalance, model training, and performance evaluation. The achieved results highlight improved classification performance through the HAN model, demonstrating its potential for enhancing software bug report prioritization.. The proposed model is evaluated on the following metrics : Precision, Recall and F1-score and it is compared with the state-of-the-art methods.

Keywords : Software Bug Reports, Bug report Prioritization, DistillBERT, Heirarichal Attention Networks (HAN), ADASYN, Deep Learning techniques, Software Reliability.

Contents

List of Figures	x
List of Tables	1
List of Acronyms	1
List of Symbols	1
1 Introduction	1
1.1 Introduction	2
1.2 Bug or Defect	2
1.3 Software Bug Report	3
1.4 Bug Report Priority	4
1.5 Organization of the Report	6
2 A Review on Bug Report Priortization	8
2.1 Review on Bug Priortization Techniques	9
3 Statement of Problem based on Identified Research Gaps	15
3.1 Research gaps	16
3.2 Research gaps	16
3.3 Objectives	17
4 Proposed Methodology	18
4.1 Proposed methodology	19
4.1.1 Working Methadology of HAN	22
4.2 Dataset	27
4.3 Evaluation Measures	28

4.4	Work Done	28
5	Results	32
5.1	Results	33
5.2	Performance Analysis	35
6	Conclusions and Future Scope	36
6.1	Conclusions	37
6.2	Future work	37
	Bibliography	38

List of Figures

1.1	Example Bug Report	5
4.1	Text preprocessing	19
4.2	Transformer Architecture	23
4.3	HAN Architecture	24
4.4	Word Level Encoding	25
4.5	Sentence Level Encoding	26

1

Introduction

1.1 Introduction

As software utilization continues to expand across various domains, software systems inevitably receive a substantial influx of bug reports. These reports are typically reviewed and categorized by individuals or automated triaging systems, which assign priorities to ensure that critical issues are promptly addressed. However, the manual assignment of priorities is a labor-intensive endeavor, prompting the need for an efficient solution. This paper proposes an automated bug report prioritization approach using deep learning models to alleviate the challenges associated with manual prioritization.

Businesses and enterprises frequently release software with imperfections due to the limitations of testing procedures. Developers rely on users to report defects using issue tracking platforms such as JIRA, GitHub Issue Tracker, Bugzilla, Mantis, and Google Code Issue Tracker. These platforms often employ unstructured natural language descriptions to capture bug details.

The act of reporting defects serves as a valuable contribution to the software maintenance process. A pivotal aspect of software maintenance involves effectively managing and addressing user-reported issues. As software systems increase in scale and complexity, the influx of problem reports grows exponentially, amplifying the intricacy of resolution efforts. Statistics reveal that during the months of January to July 2000, both Mozilla and Eclipse encountered an average of approximately 170 and 120 new bug reports, respectively.

This paper aims to streamline the prioritization process by leveraging advanced deep learning models, contributing to more efficient bug resolution and enhanced software quality.

1.2 Bug or Defect

As software applications are developed, systems regularly display unexpected behaviour and side effects that might be linked to bugs or flaws. These issues typically

arise when creating software systems. Even though some technologies can analyse and record these errors, repairing them can be a difficult and time-consuming task. Simply described, a software bug is a flaw, error, or defect in a program or system that causes expected behaviour. A "defect" is described as a flaw or malfunction in a system component that prevents it from meeting its intended requirements and calls for repair or replacement in accordance with the ISO/IEC/IEEE 24765 standard. "Bug" and "defect" are synonymous concepts. The opposite is a "bug" in a circumstance or weakness in a system that, if used or activated, could possibly cause an error in a component. In essence, multiple file problems could cause a single issue.

Repairs are actions taken to address one or more flaws, whereas debugging is a technique used to identify and fix issues. A fix frequently entails the elimination of a bug, although it can also consist of several smaller fixes to address a significant bug. The end result is the creation of "co-fixes," which may solve a number of issues in various files. Since defects may have significant effects despite being the result of minor code flaws, they can be difficult to locate and correct. There are very few times when a user cannot detect a bug in the code. As an illustration, a bug could be present in a part of code that is either never run or is only ever performed under unusual circumstances, resulting in dead code. It's possible that a number of input values or settings came together to generate this issue. A bug appears when all the requirements for a flaw or the sequence of circumstances required to bring about a flaw or imperfection in the system come together.

1.3 Software Bug Report

Developers and maintainers frequently deal with a large volume of bug reports when creating and maintaining software systems. To effectively resolve these problems, they must assess and rank the bug reports in accordance with their levels of seriousness, giving more weight to critical and severe flaws. A thorough report detailing the problems with the software and potential fixes is known as a bug report. The term "defect report" as well as "fault report," "failure report," "error report," "issue report," and "trouble report" can

1. Introduction

all be used to refer to this document. The bug report often includes crucial information including the bug's identification number, title, date of reporting and modification, severity, the developer in charge of resolving the flaw, and the bug's status, including whether it has been fixed or not. Figure 1.1 contains an illustration of a bug report.

Bug reports are necessary for keeping the software operational and play a critical part in maintaining the product's performance. Data mining techniques including clustering, classification, and information retrieval have been used in several studies over the past ten years to prioritise problem reports. These methods aid in preventing improper bug prioritisation and guarantee that critical issues are fixed quickly and effectively. Software developers and maintainers can more effectively direct their resources and efforts towards fixing serious and important flaws by analysing and prioritising bug reports, which enhances the functionality of the software as a whole.

1.4 Bug Report Priority

Bug prioritising determines the order in which issues should be repaired based on their severity, which reflects the level of risk they pose and priority level. It's crucial to start with the most serious flaws and assign the proper priorities in order to handle big software maintenance quickly. Bug reports should be given the appropriate priority throughout the prioritising process because a delay in correcting major problems may have severe repercussions. The priority level given to a problem determines how quickly it will be solved. QA analysts and software testers must therefore understand the ideas of priority in order to manage their work. Bugs can have a priority level ranging from P1, which denotes the highest priority, to P5, which denotes the lowest importance. There are five categories for bug priority levels: P1, P2, P3 P4. and P5 [2]. They priority levels are explained below:

- P1-A bug categorized as P1 or "Blocker" denotes a critical issue that requires immediate attention as it can cause the application to malfunction or become unusable.

Bug Report: #Bug-ID

Bug ID	#Bug-ID (e.g. #test_ABC, #123, #home_123)—
Tester	<i>Han Solo</i>
Date (submitted)	<i>01.11.2018</i>
Title	<i>CONTACT FORM - No confirmation message is shown</i>
Bug Description	
URL	<i>https://milleniumfalcon.rebells/contact-us</i>
Summary	The page where the contact form is on does not show any confirmation message after submitting a contact request.
Screenshot	see attached screenshots
Platform	Mac OS 10.14
Browser	Chrome 69.0.3497.100 (Official Build) (64-bit)
Administrative	
Assigned To	<i>Chewbacca</i>
Assigned At	<i>01.11.2018</i>
Priority	High
Severity	Critical

Additional Notes:

- **Step-by-step Description:**
 - a. Filled out contact form
 - b. clicked on submit
 - c. form loaded a while
 - d. contact form showed "<empty>"
- **Screenshots:**
- **Result - Should:** "Thank you for contacting us, we will get in touch shortly"
- **Result - Is:** "<empty>"

Figure 1.1: Example Bug Report

- P2-Bugs can significantly impact the fundamental functionality of one or multiple application, thus resolving them expeditiously within a few days is imperative.
- P3- The matter at hand is not urgent and can be dealt with during the regular course of activities, for instance, in the second sprint.
- P4-The defect is not of high severity and can be dealt with once the more important bugs are resolved.
- P5-Bugs that have no impact on the core functionality of the software and may or may not be addressed in the future.

1.5 Organization of the Report

The report is divided into 5 chapters. Chapter 1 consisted of Introduction. A brief description of what the remaining chapters consist of is discussed below:

- Chapter 2 mentions about the existing methods for Bug Report Prioritization.
- Chapter 3 describes the Research Gaps, the problem statement and the objectives of the proposed method.
- Chapter 4 discusses the proposed methodology and the work done.
- Chapter 5 discusses the Results
- Chapter 6 discusses the Conclusions and the Future Scope.

2

A Review on Bug Report Prioritization

In the initial phase of the research, a major focus has been given to reviewing existing state-of-the-art classification methods. Overall task can be divided into following major subtask, i.e, identifying the existing model implementations for the problem, looking for new ideas to implement on the classification task. Several methods have already been used for classification of bug reports priority, but there is a scope for improvement with the newer concepts like Hierarichal Attention Networks(HAN), and transformer model like BERT.

2.1 Review on Bug Priortization Techniques

Tan et al. [6] proposed a novel method of determining severity. The Mozilla, Eclipse, and GCC bug reports were linked to the bug repository entry on stack overflows. To determine the severity of defects, three classification algorithms-KNN, Naive Bayes, and Long Short-Term Memory (LSTM) were used. The trials findings revealed a rise in the average F-measurement of Mozilla, Eclipse, and GCC using the suggested method of 23.03 percent, 21.56%, and 20.59%.

Xiao et al. [7] for automatic bug localization, DeepLoc, an improved Convolutional Neural Networks (CNN)-based model, was proposed. CNN-inspired word embedding techniques were used by DeepLoc to replace the functionality of bug reports and source files. The experiments were conducted on 18,500 bug complaints that were extracted from five projects: Aspect J, Eclipse UI, JDT, SWT, and Tomcat between 2001 and 2014. They compared the effectiveness of DeepLoc to the four bug localization strategies BugLocator, LR+WE, HyLoc, and DeepLoc. According to the trial findings, using DeepLoc has increased Mean Average Precision (MAP) for bug localization from 10.87% to 13.1% when compared to regular CNN.

Shakripur et al. [8] employing the Time TF-IDF weighting method, the ABA-TF-IDE time-based methodology was proposed. The information was gathered from the Version Control System (VCS) software repository, where other project details and source code updates are kept. The model was trained using four machine learning algorithms:

2. A Review on Bug Report Prioritization

Support Vector Machine (SVM), Naive Bayes, Vector Space Model (VSM), and Smooth Unigram Model (SUM). The findings demonstrate that the suggested method functioned satisfactorily, with Mean Reciprocal Ranks (MRR) up to 11.8% and 8.94%.

Lamkanfi et al. [9] investigated whether it is possible to forecast the severity of a reported defect by studying its textual description using text mining methods. They used a crude Bayes algorithm on historical information gathered from the open-source community (Mozilla, Eclipse, and GNOME). They found that GNOME performed better than Mozilla and Eclipse, with precision and recall ranging from 0.70-0.85 for GNOME. Abdelmoez et al. [10] proposed a method to forecast the priority of bug reports using naive Bayes classifier. They made use of data from four different systems from the Mozilla, Eclipse, and GNOME open-source projects, which are all very large. The bug reports were ranked in order of mean time.

Qasim et al. [11] developed an automated classification method (DRONE) for forecasting bug reports' priority. They used linear regression (LR) to classify priorities and got an average F1-score of up to 29.

Alenezi and Banitaan et al. [12] used decision trees, random forests, and naive Bayes to carry out the priority prediction. They employed two feature sets: one based on the TF-weighted words of issue reports and the other on the attributes used to categorise bug reports. They employed two feature sets: one based on the TF-weighted words of issue reports and the other on the attributes used to categorise bug reports. The findings of the evaluation point to the second feature set as being more effective than the first, where naive Bayes is outperformed by decision trees and random forests.

Tian et al [13] using the nearest neighbour method to find precise bug report labels, it was possible to forecast the priority of bug reports. They used the suggested method on a bigger set of problem reports, which included over 65,000 Bugzilla complaints. Hui et al. [1] proposed a convolutional neural network based automatic approach to predict the multiclass priority for bug reports. The proposed approach significantly outperforms the state-of-the-art approaches and improves the average F1-score by more than 24%

Choudhary et al. [14] proposed an SVM-based model for priority prediction was recently created, and it prioritises Firefox crash reports in the Mozilla Socorro server according to the frequency and entropy of the crashes Table 2.1 describes few of the techniques used for bug prioritization and their limitations.

2. A Review on Bug Report Prioritization

Technique Names	Description	Limitations
CLASSIFICATION DRONE with GREY engine [15]	Labels the priority of bug reports. Contains new GREY classification engine that improves linear regression with our thresholding method to handle unbalanced bug report data.	The validation data differs from the test data, therefore the learning thresholds may not be erroneous
Information-gain and Chi-square are contrasted with Naive Bayes Multinomial (NBM) and K closest Neighbour (KNN). [16]	The TDM matrix of the dictionary of words is trained using algorithms after the most informative phrases are selected from a corpus of the matrix.	Information gained: The severe categories are less precise than the non-severe ones. Chi square: SWT component's classification of bug severity has the lowest accuracy
Context-Based Prioritization using impact analysis technique VSM [17]	Uses an issue tracking system to help issue-driven software development projects keep track of their list of issues,	The study was limited to one project and one metric, which has an influence on the reliability of the findings.
Support Vector Machine (SVM) [11]	A supervised machine learning method that classifies and carries out regression on data Utilises a kernel approach to alter data and identify the best border between two potential outputs.	Performance across all priority levels is inconsistent, and trivial reports were given P1 priority.
Logistic Regression (LOGR) [18]	With the use of Logistic Regression. which employs a number of independent factors as input to forecast the categorical dependent variable as output. it is possible to use the taxonomy of bugs to automatically classify them	The target label is more likely to forecast less accurate outcomes if it does not correlate with the problem report features.

Table 2.1: Review on Bug Prioritization

Technique Names	Description	Limitations
K Nearest Neighbor, TF-IDF [19]	The method needs stack traces from the bug data set as well as its categorical attributes in order to assess a problem's severity	It employs instance learning and is a lazy learner. Every time, the temporal complexity is $O(n)$.
Random Forest (RF) (20)	Methods for analysing and categorising code according to Clean, Smelly, Buggy, and Harmful categories, boosting programmes like Ada Boost	The algorithm can become ineffective when dealing with a lot of trees since it can be quite sluggish.
Convolutional Neural Network (CNN) [21]	An automated method to forecast the multi-class priority of bug reports,	The classifier is much slower since it is the Maxpool dataset, and processing and training the dataset on a poor GPU is quite time-consuming.

Table 2.2: Review on Bug Priortization

3

Statement of Problem based on Identified Research Gaps

3.1 Research gaps

The method of problem/defect prioritising largely relies on automated approaches that can help developers and other stakeholders process a large volume of bug reports. These systems allow for the anticipation of a lang's importance and the discovery of the bugs in the code base using specialised algorithms. In a business setting, it is useful because these techniques can cut down on time. The cost of maintenance may be reduced because ML algorithms may automatically assign the appropriate bug severity and assign the developers with the relevant abilities to fix an imur. The outcome from these preliminary steps is used as training data for ML classifiers. In order to prioritise bugs properly, a number of actions must be carried out in the proper sequence. It cannot provide effective outcomes if these first stages are not effectively carried out or if some procedures, are skipped. Few challenges faced are :-

- Managing new bugs: When bug reports are generated, it is discovered that the priorities given are inconsistent.
- Time Requirement - Some models use instance learning, which is time-consuming to run and assign priorities to.
- Bug report evaluation - Analysing the reports to decide what level of importance should be given to the bug
- Computation difficulty - The complexity of computation increases as the dataset gets bigger.
- Inconsistent findings: Testing on new datasets produces erratic results.

3.2 Research gaps

Given a bug report the aim is to develop a model which can assign priorities to bug reports which gives accurate results with good consistency and can handle new bug reports with consistency.

3.3 Objectives

The proposed method aim to fulfill the following Objectives

- To create a model that can automatically prioritise bug reports using deep learning approaches after preprocessing the data using NLP techniques.
- To compare the outcomes of the proposed method to the state-of-the-art methodology after training the suggested model on open-source bug reports from Mozilla, Eclipse, etc.

4

Proposed Methodology

4.1 Proposed methodology

By downloading bug reports from Bugzilla, a bug report dataset is constructed in the first phase. Long reports from the Eclipse JDK, Eclipse Platform, Mozilla.com, and Mozilla Firefox are created in this document. The columns in this data include Issue.id and Component. Title of the Bug, The bug's description and the priority given to the bug report. The dataset's features are utilised to assess the severity of the problem. Title, Description, and other features will be provided to the model as inputs, and the model will then predict the label, which in this case is priority. The block diagram for dataset preprocessing is shown in Fig. 4.1.

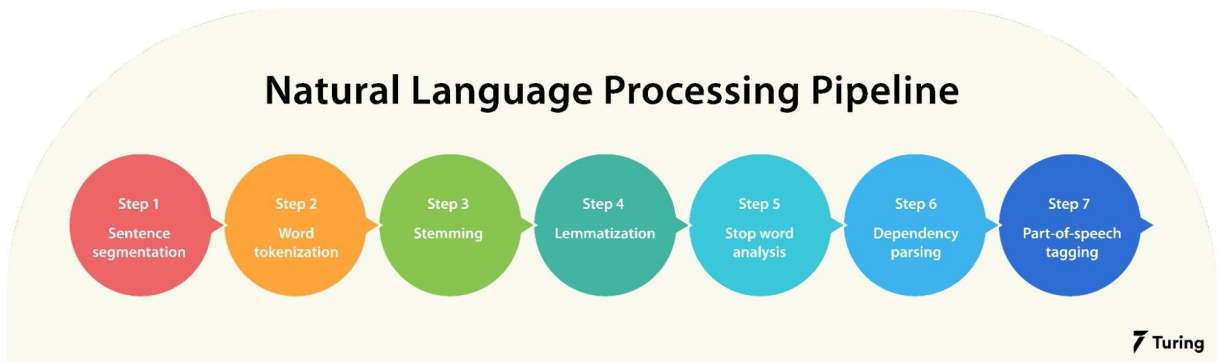


Figure 4.1: Text preprocessing

To complete this task of automatic assigning of priority to bug reports in this paper we propose a novel prototype in which data is pre-preprocessed using techniques, such as Tokenization, Spell correct, Stop-word removal, Lemmatization and then is data is passed to DistilBert Tokenizer. The Distillert tokenizer use WordPiece algorithm to break down input text into individual tokens, which are then converted into numerical representations that can be processed by machine learning model. These numerical representations are passed to DistilBert model for training Fig 4.2 shows the block diagram of the proposed model.

DistilBert is a smaller, distilled version of the popular BERT (Bidirectional Encoder Representations from Transformers) model. DistillBert tokenizer will be used in the text preprocseing. The key innovation behind DistilBert is that it uses a "distillation"

4. Proposed Methodology

technique to compress the larger BERT model into a smaller version. During this process, the model's architecture is simplified, and some of its parameters are removed. TDistilBert uses a WordPiece tokenizer to break input text into tokens, which are then represented as vectors, Fig 4.1 shows the diagram of example of Bert tokenizer.

Also for Word embedding GloVe is used. GloVe (Global Vectors for Word Representation) is a word embedding technique. It utilizes word co-occurrence statistics from text data to create word vectors that encode semantic relationships. Key aspects of GloVe include:

- **Word Co-Occurrence Matrix:** GloVe constructs a matrix capturing word co-occurrence frequencies in a text corpus context window.
- **Word Vectorization:** GloVe transforms the co-occurrence matrix into word vectors through optimization, aiming to predict co-occurrence probabilities.
- **Loss Function and Optimization:** A loss function measures the difference between predicted and actual co-occurrence probabilities. An iterative process optimizes word vectors.
- **Final Word Vectors:** The learned vectors represent words, preserving semantic relationships and placing similar words closer in vector space.

Thus, this paper proposes the use of DistilBert as tokenizer for text preprocessing and Hierarchical Attention Network, which is attention based mechanism which when trained on dataset possibly provide better data than existing methods. Fig 4.1 shows the DistilBert architecture.

Hierarchical Attention Networks (HAN) uses stacked recurrent neural networks on word level followed by attention model to extract such words that are important to the meaning of the sentence and aggregate the representation of those informative words to form a sentence vector. Then the same procedure applied to the derived sentence vectors which then generate a vector who conceives the meaning of the given document and that vector

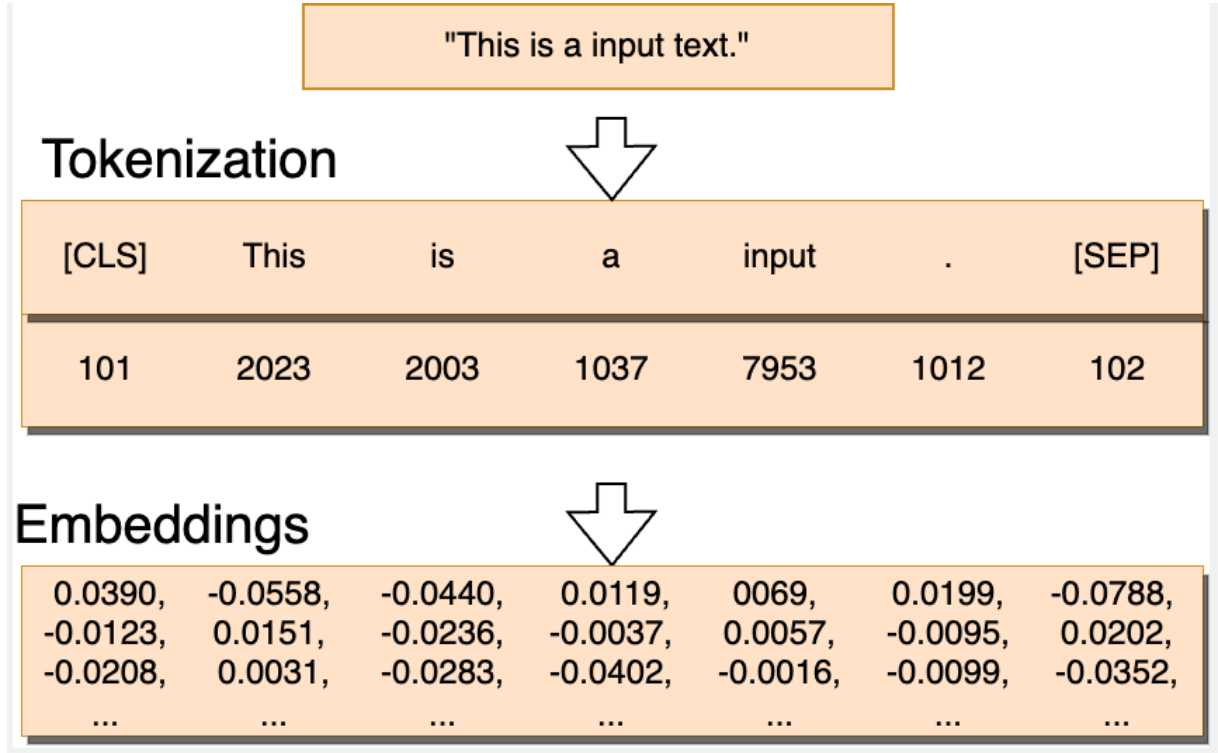


Figure 4.1: DistillBert-Tokenizer

can be passed further for text classification.

The Hierarchical attention is based on the concept that : "Words form sentences, and sentences form documents." To capture sentence and document meanings, an attention model is used, giving importance to key words. This model involves Bidirectional RNN and Attention networks. Bidirectional RNN interprets word sequences, yielding individual word vectors. Attention network assigns weights to these vectors via a neural network, aggregating them into a sentence vector through weighted summation. This process extends to sentence and document levels, forming a hierarchical attention network. The vectors undergo a shallow neural network to determine weights, and their weighted sum captures combined meaning.

4. Proposed Methodology

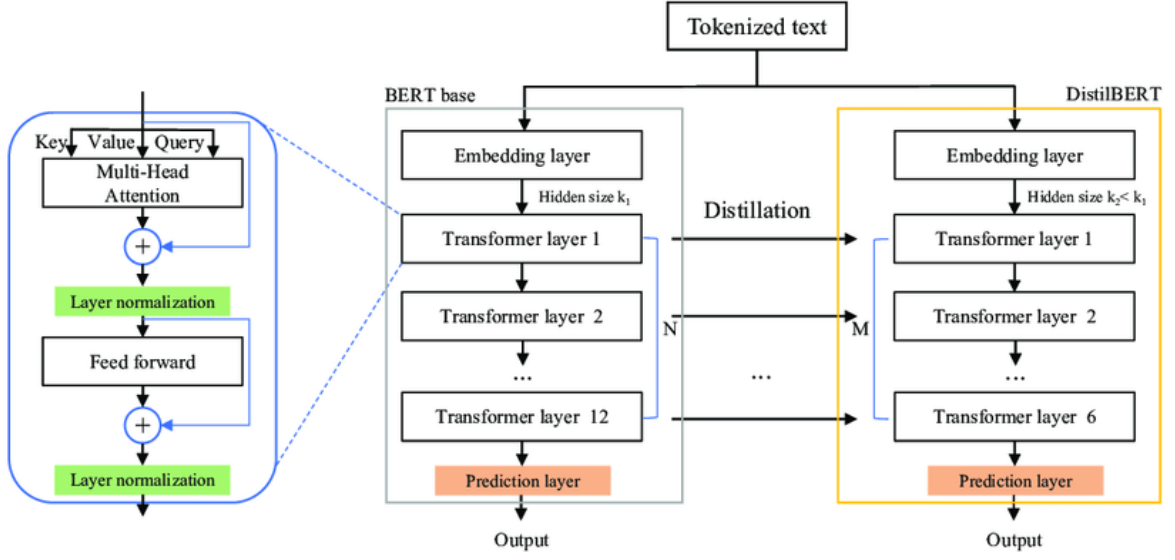


Figure 4.1: DistillBert Architecture

4.1.1 Working Methodology of HAN

First, the network considers the hierarchical structure of documents by constructing a document representation by building representations of sentences and then aggregating those into a document representation. 'Sentence representations' are built by first encoding the word of a sentence and then applying the attention mechanism on them resulting in a *sentence vector*. 'Document representation' is built in the same way, however, it only receives the sentence vector of each sentence of the document as input. The algorithm is applied twice as: First on word level and afterwards on sentence level. The model consists of

- the encoder which returns relevant contexts, and
- the attention mechanism, which computes importance weights of these contexts as one vector.

Word Level

Referring Figure 4.4

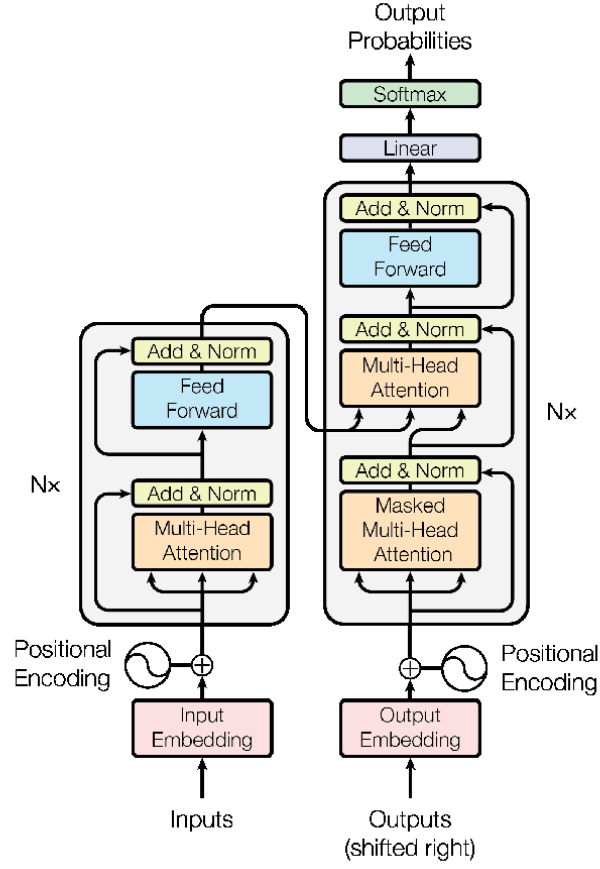


Figure 4.2: Transformer Architecture

$$x_{it} = W_e w_{it}, \quad t \in [1, T]$$

- As input the structured tokens ****wit **** which represent the word i per sentence t are provided. We do not keep and process all words in a sentence.
- Since the model is not able to process plain text of the data type string, the tokens run through an Embedding layer which ‘assigns’ multidimensional vectors ****We **** ****wit **** to each token. In this way, words are represented numerically as ****xit **** as a projection of the word in a continuous vector space.
- There are several embedding algorithms: the most popular ones are word2vec and GloVe. Pretrained Glove embeddings have been used for this methods.

4. Proposed Methodology

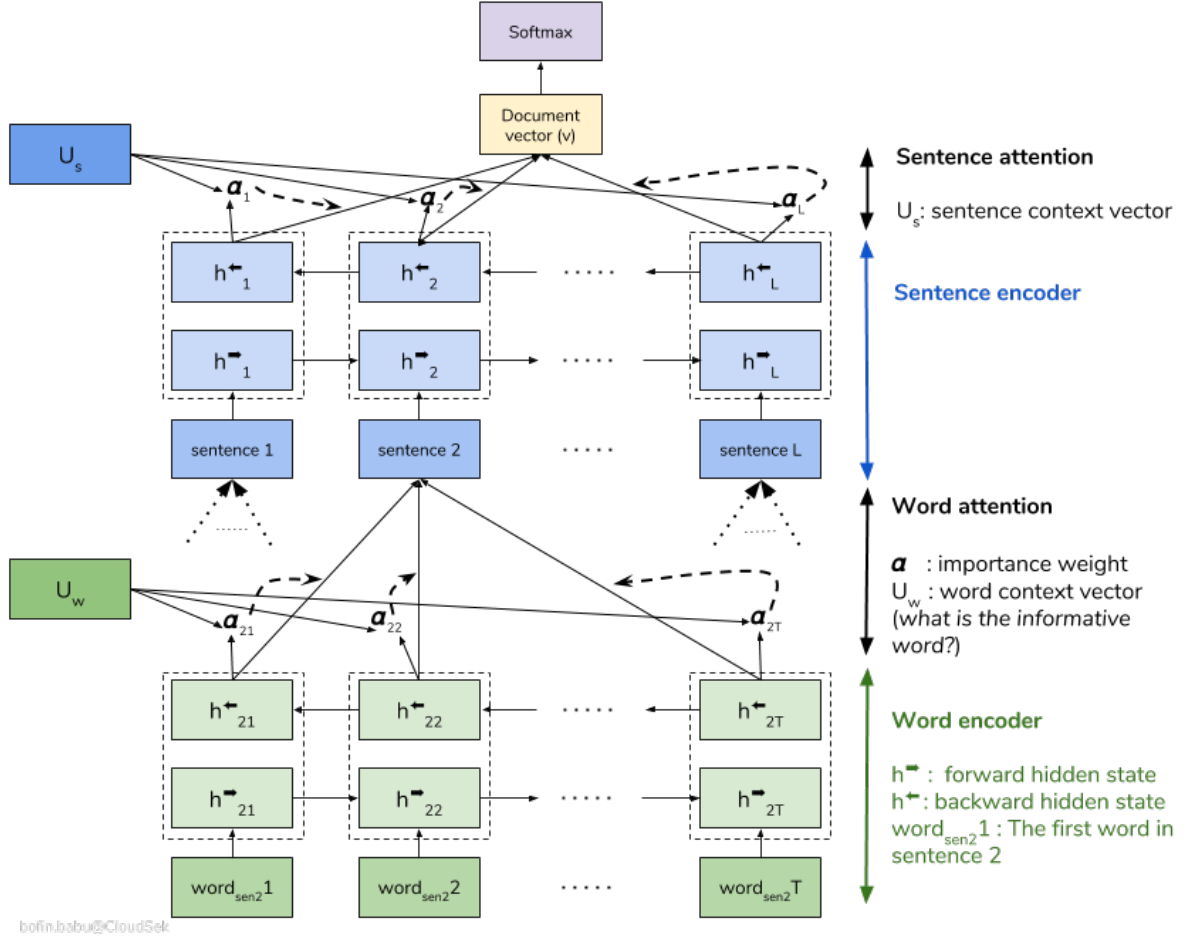


Figure 4.3: HAN Architecture

Word Encoder

Given a sentence with words w_{it} , $t \in [0, T]$, first embed the words to vectors through an embedding matrix W_e , where $x_{ij} = W_e w_{ij}$. We use a bidirectional GRU (Bahdanau et al., 2014) to obtain annotations of words by summarizing information from both directions for words, incorporating contextual information in the annotation. The bidirectional GRU contains the forward GRU, vec_f , which reads the sentence S_i from w to w_{iT} , and a backward GRU, vec_b , which reads from w_{iT} to w_{i1} .

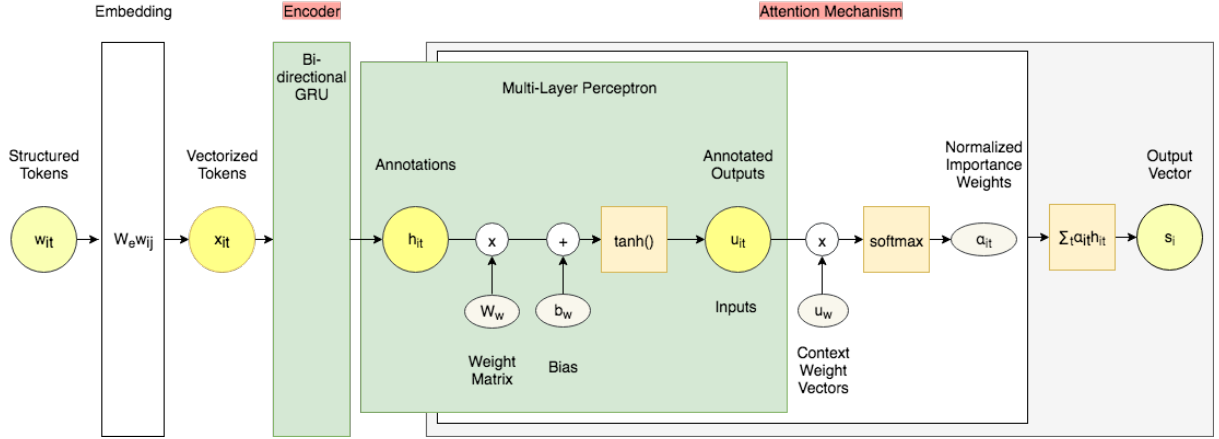


Figure 4.4: Word Level Encoding

$$\begin{aligned}
 x_{it} &= W_e w_{it}, & t &\in [1, T] \\
 \text{vec}_{h_{it}} &= \text{vec GRU}(x_{it}), & t &\in [1, T] \\
 \bar{h}_{it} &= \overline{\text{GRU}}(x_{it}), & t &\in [T, 1]
 \end{aligned}$$

annotation for a given word is obtained w_{it} by concatenating the forward hidden state $\text{vec}_{h_{it}}$ and the backward hidden state \bar{h}_{it} , i.e., $h_{it} = [\text{vec}_{h_{it}}, \bar{h}_{it}]$, which summarizes the information of the whole sentence centered around w_{it} .

Word Attention

- The annotations \mathbf{h}_{i1} build the base for the attention mechanism, which starts with another hidden layer, a one-layer Multilayer Perceptron. The goal is to let the model learn through training with randomly initialized weights and biases. Those 'improved' annotations are then represented by \mathbf{u}_{it} . Furthermore, this layer ensures that the network does not falter with a tanh function. This function corrects input values to be between -1 and 1 and also maps zeros to near-zero.

$$\mathbf{u}_{it} = \tanh(\mathbf{W}_w \mathbf{h}_{it} + \mathbf{b}_w)$$

4. Proposed Methodology

- Our new annotations are again multiplied with a trainable context vector \mathbf{u}_w and normalized to an importance weight per word \mathbf{C}_L , by a softmax function. The word context vector \mathbf{u} is randomly initialized and jointly learned during the training process.

$$\alpha_{it} = \frac{\exp(\mathbf{u}_{it}^\top \mathbf{u}_w)}{\sum_t \exp(\mathbf{u}_{it}^\top \mathbf{u}_w)}$$

- The sum of these importance weights concatenated with the previously calculated context annotations is called the sentence vector.

$$\mathbf{s}_i = \sum_t \alpha_{it} \mathbf{h}_{it}$$

Sentence Level

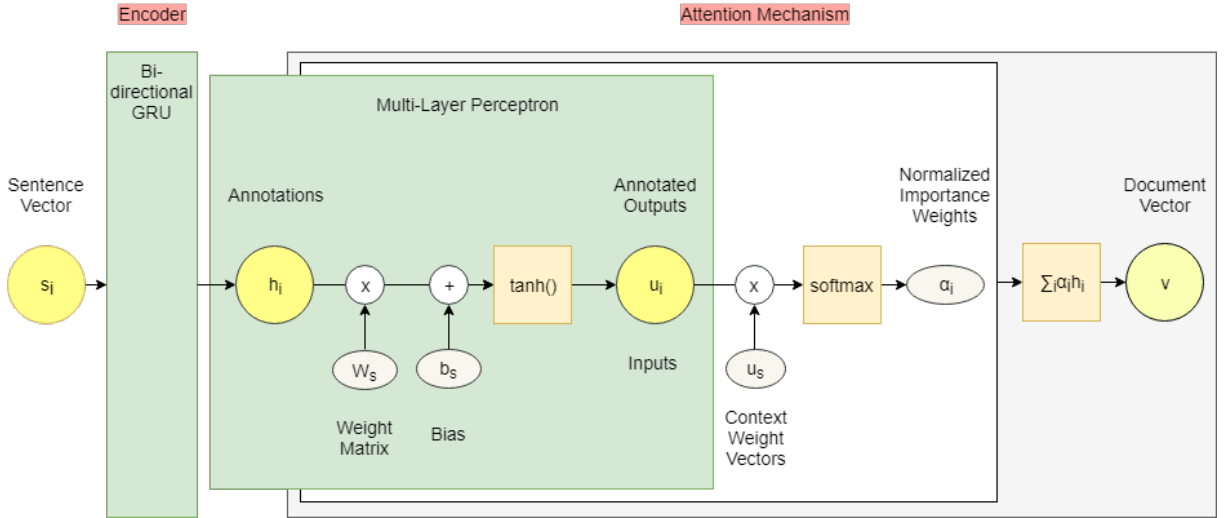


Figure 4.5: Sentence Level Encoding

Then the whole network is run on sentence level with basically the same procedure as on word level, but now we focus on the sentence i . There is no embedding layer as we already get sentence vectors \mathbf{s}_i from word level as input.

Sentence Encoder

Contexts of sentences are summarized with a bidirectional GRU by going through the document forwards and backwards.

$$\vec{h} = \text{GRU}(\mathbf{s}), \quad t \in [1, L]$$

$$\overleftarrow{h} = \text{GRU}(\mathbf{s}), \quad t \in [L, 1]$$

$$h_1 = [\vec{h}_i, \overleftarrow{h}_i]$$

Sentence Attention

- Trainable weights and biases are again randomly initialized and jointly learned during the training process. The final output is a document vector \mathbf{v} which can be used as features for document classification.

$$u_i = \tanh(W_{sh}\mathbf{i} + b_s)$$

$$\alpha_i = \frac{\exp(u_i^\top u_s)}{\sum_t \exp(u_t^\top u_s)}$$

$$v = \sum a_i h_i$$

4.2 Dataset

A bug report dataset is created by downloading bug reports from Bogzilla. In this paper we have downloaded the dataset from a Github repository which has a collection of all these bug reports and is freely available for research purposes. In this paper bog reports from eclipse jdt, eclipse-platform, mozilla.core, mozilla firefox is created. This dataset has several columns such as Issue.id, Component, Title of the bog, Description of the bug and Priority assigned to the bug report. These features of the dataset are used to

4. Proposed Methodology

determine the priority of the bug. Here, priority will be the label and Title Description, etc. will be the features which will be given as input to the model and the model will predict the label L_e , the priority. This dataset consists of bug reports from the open-source softwares mentioned above. Four csv files are generated for each of the above mentioned open-source softwares, eclipse.jdt, eclipse platform, mozilla.core. Mozilla firefox.

4.3 Evaluation Measures

The proposed model is evaluated on the following metrics: Precision, Recall and F1-score and the model is compared with the state-of-the-art methods:

- Precision: In terms of positive observations, precision is the proportion of accurately anticipated observations to all expected positive observations.
- Recall Recall is defined as the proportion of accurately predicted positive observations to all of the actual class observations.
- FI-score: The FI-Score is a measurement that combines recall and accuracy. The weighted average of Precision and Recall is the FI Score. Therefore, both false positives and false negatives are included while calculating this score.

4.4 Work Done

We tested the proposed model on eclipse platform dataset on one feature which is the Description of the bug feature. Dataset was created by downloading several bug reports of eclipse from Bugzilla and combining. The dataset has the following columns: Issue.id, Title, Description, Status, Resolution and Created-time and the label Priority. The data was preprocessed and then vector representations were obtained using Distilbert Tokenizer. The data was then oversampled using ADASYN. Glove Word embedding was used for the pretrained numerical values. The glove numerical representations were passed on the HAN model and the model was trained.

Text Preprocessing

The textual data underwent a series of preprocessing steps to enhance its quality and relevance for model training. The following preprocessing steps were applied to each text sample:

- **Lowercasing:** All text was converted to lowercase to ensure consistency in word representation.
- **Punctuation Removal:** Punctuation marks and special characters were removed from the text.
- **Stopword Removal:** Common stopwords from the English language were eliminated from the text to focus on meaningful content.
- **Stemming:** Words were stemmed using the Porter stemming algorithm to reduce them to their root forms.
- **Lemmatization:** Lemmatization was applied to further normalize words by converting them to their base or dictionary form.

These preprocessing steps aimed to create a cleaner and more standardized textual dataset for model training.

GloVe Word Embeddings

GloVe (Global Vectors for Word Representation) embeddings were employed to enhance the representation of words in the textual data. These pre-trained word vectors capture semantic relationships between words and help improve the model's understanding of the textual content. The GloVe embeddings were loaded from a pre-trained file (`glove.6B.100d.txt`) and utilized for initializing the embedding layer of the HAN model.

4. Proposed Methodology

ADASYN Oversampling

To address class imbalance in the dataset, the ADASYN (Adaptive Synthetic Sampling) technique was employed. This method oversamples the minority class by generating synthetic samples that are similar to existing samples. ADASYN helps improve the model's ability to learn from underrepresented classes and avoid bias towards majority classes, thereby enhancing the overall performance of the model on all classes.

Hierarchical Attention Network (HAN) Model

The Hierarchical Attention Network (HAN) model was constructed to effectively capture the hierarchical structure of textual data. The model consists of two levels of attention mechanisms: word-level attention and sentence-level attention. The HAN model processes text in a hierarchical manner, allowing it to focus on important words and sentences while understanding the context and relationships within the text.

The model architecture consists of the following layers:

- **Embedding Layer:** Initialized with pre-trained GloVe embeddings to represent words in a dense vector space.
- **Bidirectional GRU Layers:** These layers encode both forward and backward contextual information for words and sentences.
- **Attention Mechanisms:** Word-level and sentence-level attention mechanisms highlight relevant words and sentences within the text.
- **Dense Layers:** Fully connected dense layers for feature extraction and classification.

Training and Evaluation

The model was trained on the preprocessed and oversampled dataset using a batch size of 64 and a maximum of 10 epochs. The training process was monitored using various performance metrics, including accuracy, loss, precision, recall, and F1-score. Early stopping was employed to prevent overfitting and optimize training efficiency.

Plots and Visualizations

Several plots and visualizations were generated to analyze and interpret the model's performance:

- **Accuracy and Loss Plots:** Plots depicting the training and validation accuracy, as well as the loss, across the training epochs.
- **Confusion Matrix:** A heatmap representation of the confusion matrix, displaying the distribution of predicted labels against true labels for each class.
- **Precision-Recall-F1 Curves:** Precision, recall, and F1-score were calculated and presented for each class, providing insights into the model's class-specific performance.

Both undersampling and oversampling were used in the model's training process. In the situation of undersampling, the performamodel did not perform very well. The model performed well after oversampling, though. As a result, the dataset was subjected to the oversampling approach known as ADASYN. The model was tested on a test dataset after the training phase. The loss metric utilised was sparse categorical cross-entropy, and the optimizer employed was the Adam optimizer. categorical cross-entropy that is sparse. For multi-class classification, sparse categorical cross-entropy is utilised. A development of the stochastic gradient descent method for updating network weights during training is the optimisation algorithm known as Adam. Adam changes the learning rate for each network weight independently, as contrast to (SGD) Stochastic Gradient Descent, which utilises a single learning rate for the entirety of the training process. 128-piece batches were used to train the model. Ten epochs were used to train the model. The model's outcomes were contrasted with those of previously employed methods.

5

Results

5.1 Results

The evaluation metric used for this purpose is F1-score and it is evaluated on the priority levels. F1-score for every priority level is calculated and compared. F1-core is defined as the harmonic mean of precision and recall.

Precision (Positive Predictive Value)

Precision measures the proportion of correctly predicted positive instances (true positives) out of all instances predicted as positive (true positives + false positives).

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Recall (Sensitivity, True Positive Rate)

Recall measures the proportion of correctly predicted positive instances (true positives) out of all actual positive instances (true positives + false negatives).

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

F1-Score

The F1-score is the harmonic mean of precision and recall. It provides a balance between precision and recall, making it useful for evaluating models when there is an imbalance between the classes.

$$\text{F1-Score} = \frac{2 \cdot (\text{Precision} \cdot \text{Recall})}{\text{Precision} + \text{Recall}}$$

Here for every priority level Lc. P1, P2, P3, P4 and P5 Precision, Recall and El-score is calculated. The proposed model is evaluated on priority level and it is compared with existing approaches. For the priority levels P1, P2, P3, P4 and P5 the results obtained are as follows:

5. Results

- Precision values obtained for the proposed model: 54.84%, 28.57%, 38.71%, 37.14%, and 38.33% respectively.
- Recall values obtained for the proposed model: 32.69%, 33.33%, 23.08%, 53.06%, and 47.92% respectively.
- F1-score obtained for the proposed model: 40.96%, 30.77%, 28.92%, 43.70%, and 42.59% respectively.

In the result proposed model is compared with 4 base models: Logistic Regression, SVM, Naive Bayes and CNN. It is homogenous comparison as the dataset and preprocessing steps etc., used are same. The results obtained by the proposed method and its comparison with other base models implemented (homogenous comparison) is displayed in Tables 5.1, 5.2 and 5.3. Table 5.1 shows Precision, Table 5.2 shoes Recall and Table 5.3 shoes F1-Score.

Table 5.1: Results of the present HAN model in terms of Precision

Techniques	RESULTS (Precision)				
	P1	P2	P3	P4	P5
Proposed Method	54.84%	28.57%	38.71%	37.14%	38.33%
CNN	37.16%	33.51%	43.80%	39.71%	27.39%
SVM	35.41%	28.31%	40.38%	37.26%	25.93%
Random Forest	34.23%	30.29%	38.56%	35.09%	24.64%
Logistic Regression	32.17%	25.50%	35.68%	33.11%	22.93%

Table 5.2: Results of the presented HAN model in terms of Recall

Techniques	RESULTS (Recall)				
	P1	P2	P3	P4	P5
Proposed Method	32.69%	43.33%	23.08%	53.06%	47.92%
CNN	40.67%	21.86%	44.15%	39.77%	30.43%
SVM	36.92%	21.13%	40.54%	38.66%	27.72%
Random Forest	36.52%	20.14%	43.50%	39.59%	27.13%
Logistic Regression	33.25%	20.19%	38.40%	27.72%	25.36%

Table 5.3: Results of the presented HAN model in terms of F1-score

Techniques	RESULTS (F1-score)				
	P1	P2	P3	P4	P5
Proposed Method	40.96%	30.77%	28.92%	43.70%	42.59%
CNN	33.83%	26.46%	23.54%	39.74%	28.83%
SVM	37.02%	23.54%	31.22%	38.39%	26.52%
Random Forest	35.52%	24.89%	29.53%	36.79%	26.09%
Logistic Regression	32.70%	22.54%	26.99%	33.85%	24.08%

5.2 Performance Analysis

From the obtained results, we can observe that the F1-score obtained for priority P1. P2. P4 and P5 is better than that of the existing methods. When compared with base models (homogenous comparison) which were trained on the same dataset as that of the proposed model and also same pre-processing steps were applied, it can be seen that the performance of the proposed model is somewhat better than that of the existing models. Also, it can be inferred that when the proposed model will be trained on larger datasets its performance will significantly improve.

6

Conclusions and Future Scope

6.1 Conclusions

A review of the research work had been finished, as can be seen in chapter 2. Different research and studies have various characteristics, advantages, disadvantages, and most significantly, promise in the future. It is possible to infer from reading the review that bug reports can be automatically prioritized with the aid of machine learning and deep learning techniques, as well as their combination with NLP approaches to produce superior outcomes. While it has been found via the study that precise prioritization of problems may be very beneficial for the overall functionality of the software as well as for the programmers or maintainers, which will boost the overall effectiveness of the task

6.2 Future work

The results above show that the suggested technique performs admirably when trained on the dataset. By using better preprocessing methods, this outcome can be enhanced: To enhance the outcomes, the suggested model can also be trained on bigger datasets and for more epochs. Improved accuracy can be achieved by fine-tuning the model. The findings can be improved by enhancing the overall model, which will also improve the hyperparameter of the HAN model. Other machine learning and deep learning approaches can be investigated for this problem in additional study and contrasted with the current methods and the suggested ways.

Bibliography

