

Project Report Final

Team Members - Javed Habib (12140830), Mohd. Adil (12141080)

Aim of the project

When open-source projects receive contributions and bug reports, they often involve managing a significant number of issues. Labeling issues correctly is essential for efficient issue tracking and resolution. However, it can be a time-consuming task. The aim of this project is to develop a machine learning model that can automatically predict appropriate labels for GitHub issues based on their content.

Progress

About the dataset

It contains 150000 rows with only 2 labels which only classify each row as a bug and not a bug. It contains a body and a title, both of which have mentions and links. Hence appropriate cleaning techniques were used. We used regular expressions and then standard tokenizers to generate embeddings.

		title	body	label	full_text	distilbert-base-nli-mean-tokens_embeddings	albert-base-v1_embeddings	roberta-base-nli-stsb-mean-tokens_embeddings
0	y-zoom piano roll	a y-zoom on the piano roll would be useful.	1	title : y-zoom piano roll body : a y-zoom on t...	[-0.7301139831542969, -0.047545962035655975, 0...	[0.0026986433658748865, 0.12280294299125671, -...	[0.42477285861968994, -0.09550482034683228, -0...	
1	buggy behavior in selection	! screenshot from 2016-02-23 21 27 40 https://...	0	title : buggy behavior in selection body : ! S...	[-0.7117552757263184, 0.05117011442780495, 0.3...	[0.21210505068302155, -0.3419630825519562, -0....	[0.46860456466674805, -0.4723658561706543, 0.3...	
2	auto update feature	hi,\r \r great job so far, @saenzramiro ! : \r...	1	title : auto update feature body : hi,\r \r gr...	[-0.545153796672821, 0.023795893415808678, 0.9...	[0.07694720476865768, -0.23390068113803864, -0...	[0.7905717492103577, -0.6955427527427673, 0.66...	
3	filter out noisy endpoints in logs	i think we should stop logging requests to:\r ...	1	title : filter out noisy endpoints in logs bod...	[-0.5753872990608215, 0.24625295400619507, 0.5...	[0.14907507598400116, -0.36269834637641907, -0...	[0.34696999192237854, 0.3869624435901642, 0.26...	
4	enable pid on / pid off alarm actions for ardu...	expected behavior\r alarm actions pid on and p...	0	title : enable pid on / pid off alarm actions ...	[-0.600375771522522, -0.21507282555103302, 0.5...	[-0.005618194118142128, -0.010725936852395535, ...	[0.36912062764167786, 0.15820196270942688, 0.1...	

In [9]:

```
1 def cleaner(text):
2     text = re.sub(r"(@\[A-Za-z0-9]+)|(^0-9A-Za-z \t)|(\w+:\//\S+)|^rt|http.+?", "", text)
3     return text
4
5 def standardize_accented_chars(text):
6     return unicodedata.normalize('NFKD', text).encode('ascii', 'ignore').decode('utf-8', 'ignore')
7     print(standardize_accented_chars(df['body'].loc[2]))
```

hi,\r \r great job so far, @saenzramiro ! : \r \r an auto update feature would be nice to have.\r or alternatively a menu button to check for the latest version manually.

In [11]:

```
1 def clean_text(text):
2     text = re.sub(r'\r', ' ', text)
3     text = re.sub(r'\n', ' ', text)
4     text = re.sub(r'\n', ' ', text)
5     text = re.sub(r'\s+', ' ', text).strip()
6     text = re.sub(r'http\S+|www\S+|https\S+', '', text, flags=re.MULTILINE)
7     text = re.sub(r'@[ \w]+', '', text)
8     text = re.sub(r'^A-Za-z0-9\s.,!?', '', text)
9     return text
10
11 raw_text = "hi,\r\r \r great job so far, @saenzramiro ! : \r\r \r an auto update feature would be nice to have.
12 cleaned_text = clean_text(raw_text)
13 print(cleaned_text)
```

hi, great job so far, ! an auto update feature would be nice to have. or alternatively a menu button to check fo
r the latest version manually.

Since the size of large, we had to spend a lot of time in Phase 1 making only the embeddings. We have used 3 models to generate embeddings -

1. distilbert-base-nli-mean-tokens_embeddings (dim - 768)
2. albert-base-v1_embeddings (dim - 768)
3. roberta-base-nli-stsb-mean-tokens_embeddings (dim - 768)

```

1 model = SentenceTransformer('distilbert-base-nli-mean-tokens')
2 def generate_embeddings(text):
3     embeddings = model.encode(text, convert_to_tensor=False)
4     return embeddings.tolist()
5 # generate_embeddings('I am batman')
6 df_train['distilbert-base-nli-mean-tokens_embeddings'] = df_train['full_text'].progress_map(generate_embeddings)
7
8 model = SentenceTransformer('albert-base-v1')
9 def generate_embeddings(text):
10     embeddings = model.encode(text, convert_to_tensor=False)
11     return embeddings.tolist()
12 df_train['albert-base-v1_embeddings'] = df_train['full_text'].progress_map(generate_embeddings)
13
14 model = SentenceTransformer('roberta-base-nli-stsb-mean-tokens')
15 def generate_embeddings(text):
16     embeddings = model.encode(text, convert_to_tensor=False)
17     return embeddings.tolist()
18 df_train['roberta-base-nli-stsb-mean-tokens_embeddings'] = df_train['full_text'].progress_map(generate_embeddings)
19
20
21 df_train.head()

```

We then tried to cluster these points as it is but the **curse of dimensionality** hit us as it was taking a lot of time and my system kept crashing. Hence we used **PCA** to reduce our vector to 50 dimensions. The dataset size had jumped from **200 MBs** to **8 GBs** after the generation of embeddings etc. We have tried a lot of clustering methods but only the following gave us good results -

1. DB-SCAN (gave very bad results)
2. Agglomerative Hierarchical Clustering
3. Spectral Clustering
4. K means

Then I manually the clusters by sampling and identified the following potential new classes -

1. Features
2. Bugs
3. Auth
4. Maintenance
5. Security
6. Testing

Experiments

We tried to lable our data through various LLMs by ujsing prompt engineering and **langchain** but it was giving very poor results and took a lot of time to process.

Labelling using prompts

```
In [40]: 1 from langchain.llms import LlamaCpp
2 from langchain.document_loaders import PyPDFLoader
3 from langchain.embeddings import LlamaCppEmbeddings
4 from langchain.prompts import PromptTemplate
5 from langchain.chains import LLMChain
6 from langchain.document_loaders import TextLoader
7 from langchain.text_splitter import CharacterTextSplitter
8 from langchain.vectorstores import Chroma
9 from langchain.text_splitter import RecursiveCharacterTextSplitter
10 import os
11 from langchain.document_loaders import PyPDFLoader
12 from langchain.callbacks.manager import CallbackManager
13 from langchain.callbacks.streaming_stdout import StreamingStdOutCallbackHandler
14 from langchain.llms import LlamaCpp
15 from langchain.callbacks.streaming_stdout import StreamingStdOutCallbackHandler
16 from langchain.chains import LLMChain
17 from langchain.llms import GPT4All
18 from langchain.prompts import PromptTemplate
19 import warnings
20 warnings.filterwarnings('ignore')
```

```
In [72]: 1 df['text_full'] = df['cleaned_title'] + ' ' + df['cleaned_body']
```

```
In [75]: 1 template = """
2 Issue text : {text}
3 Question: {question}
4
5 Answer: Label : """
6
7 prompt = PromptTemplate(template=template, input_variables=['doctext',"question"])
8 callback_manager = CallbackManager([StreamingStdOutCallbackHandler()])
9 # Make sure the model path is correct for your system!
10 llm = LlamaCpp(
11     model_path="..\\LRNLP\\llama-2-7b-chat.Q4_K_M.gguf",
12     temperature=0.75,
13     max_tokens=200000000,
14     top_p=1,
15     callback_manager=callback_manager,
16     verbose=True, # Verbose is required to pass to the callback manager
17 )
```

```
In [76]: 1 text = """
2 title : tighten up cli parsing body : i have a multilanguage project, so i want to organize my tests inside a m
3 """
4 question = """Please assign only one label from teh following labels
5 1. Features
6 2. Bugs
7 3. Auth
8 4. Maintenance
9 5. Security
10 6. Testing
11 make sure not to differ from these classes. The output should only be one word"""
12 prompt = f"""
13 Issue text : {text}
14 Question: {question}
15
16 Answer: Label :
17 """
18
19 llm(prompt)
```

Testing

```
llama_print_timings:      load time =      204.29 ms
llama_print_timings:      sample time =       0.51 ms /      3 runs (   0.17 ms per token, 5917.16 tokens per s
econd)
llama_print_timings: prompt eval time =    3645.50 ms /    165 tokens (   22.09 ms per token,   45.26 tokens per s
econd)
llama_print_timings:      eval time =     112.87 ms /      2 runs (   56.43 ms per token,   17.72 tokens per s
econd)
llama_print_timings:      total time =     3785.41 ms
```

Out[76]: 'Testing'

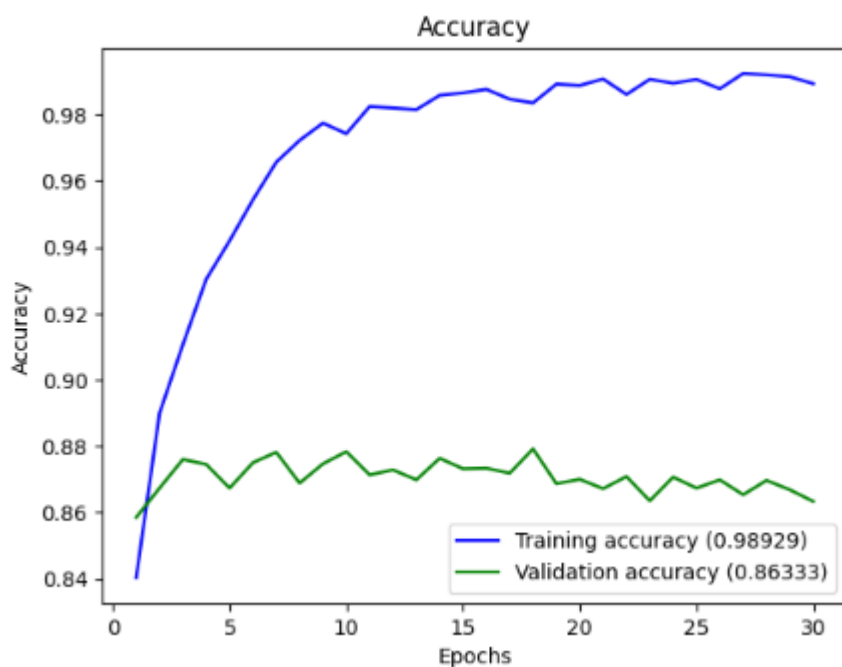
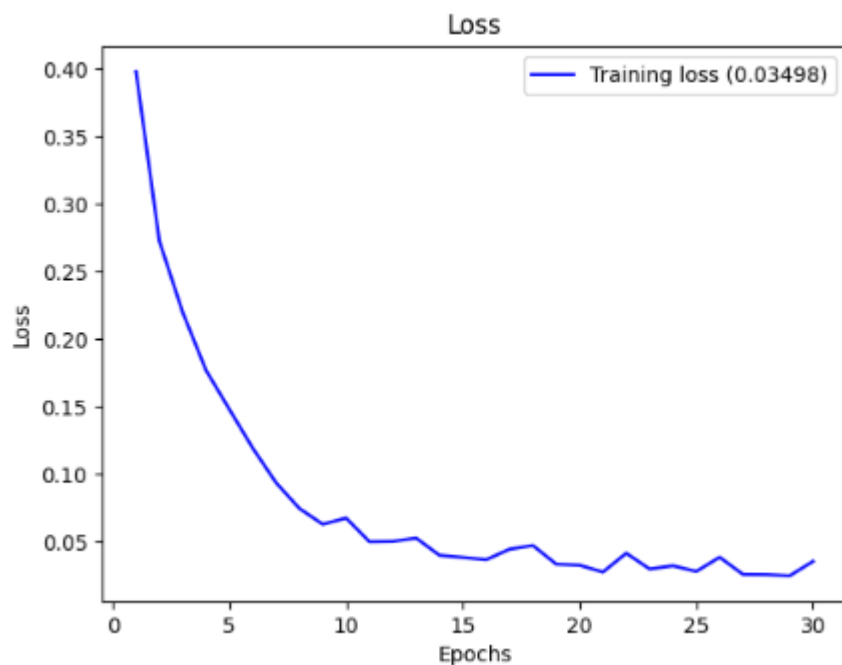
It is giving good labels...so now we can give it some data annotation task and proceed with semi supervised learning.

```
In [84]: 1 subset = df.sample(100, random_state=69)
2 subset.head()
```

Initially it was doing well but as we went on and on, the performance degraded and hence we got the following results.

```
In [87]: 1 llm_annotations
         2 'Bugs',
         3 'Security\n\nExplanation: The given text mentions the idea of extracting data generation code into an independent library, which is a common practice in software development to reuse and maintain code. The issue being discussed is labeled as "Security", indicating that it pertains to issues related to security in the context of software development.',
         4 'The GitHub issue has the label "Features".',
         5 'Bugs',
         6 'The label of the given GitHub issue is "Bugs".',
         7 'The given GitHub issue has the label "Bugs".',
         8 '"Bugs"',
         9 'Security',
        10 '035 localid is the value of tag 035 and code a when code 9 value is inspire example marcdatafield tag 035 ind1 ind2 marcsbfield code 9 inspiremarcsbfield m',
        11 'Bugs',
        12 'Security\n\nExplanation: The given GitHub issue mentions that the job is running on container-based infrastructure, which does not allow use of sudo, setuid, and setgid executables. This indicates that the issue is related to security, hence the label is "Security".',
        13 'Bugs\n\nExplanation: The given text mentions "add measurement form ability to add measurement for each type of measurement" which is a feature request. Therefore, the label for this GitHub issue would be "Bugs".',
        14 'Bugs'
```

```
accuracy: 0.8715
Epoch 12/30
750/750 [=====] - 4s 5ms/step - loss: 0.0498 - accuracy: 0.9820 - val_loss: 0.6029 - val_
accuracy: 0.8728
Epoch 13/30
750/750 [=====] - 4s 5ms/step - loss: 0.0522 - accuracy: 0.9815 - val_loss: 0.6457 - val_
accuracy: 0.8698
Epoch 14/30
750/750 [=====] - 4s 5ms/step - loss: 0.0395 - accuracy: 0.9858 - val_loss: 0.5885 - val_
accuracy: 0.8763
Epoch 15/30
750/750 [=====] - 5s 7ms/step - loss: 0.0379 - accuracy: 0.9865 - val_loss: 0.6870 - val_
accuracy: 0.8732
Epoch 16/30
750/750 [=====] - 4s 5ms/step - loss: 0.0362 - accuracy: 0.9876 - val_loss: 0.7413 - val_
accuracy: 0.8733
Epoch 17/30
750/750 [=====] - 4s 5ms/step - loss: 0.0440 - accuracy: 0.9846 - val_loss: 0.7848 - val_
accuracy: 0.8718
Epoch 18/30
750/750 [=====] - 4s 5ms/step - loss: 0.0468 - accuracy: 0.9835 - val_loss: 0.7439 - val_
accuracy: 0.8792
Epoch 19/30
750/750 [=====] - 4s 5ms/step - loss: 0.0329 - accuracy: 0.9892 - val_loss: 0.7923 - val_
accuracy: 0.8687
Epoch 20/30
750/750 [=====] - 4s 5ms/step - loss: 0.0321 - accuracy: 0.9887 - val_loss: 0.7822 - val_
accuracy: 0.8700
Epoch 21/30
750/750 [=====] - 4s 5ms/step - loss: 0.0271 - accuracy: 0.9907 - val_loss: 0.8814 - val_
accuracy: 0.8672
Epoch 22/30
750/750 [=====] - 4s 5ms/step - loss: 0.0410 - accuracy: 0.9860 - val_loss: 0.8969 - val_
accuracy: 0.8708
Epoch 23/30
750/750 [=====] - 4s 5ms/step - loss: 0.0293 - accuracy: 0.9906 - val_loss: 1.0043 - val_
accuracy: 0.8635
Epoch 24/30
750/750 [=====] - 5s 7ms/step - loss: 0.0317 - accuracy: 0.9894 - val_loss: 0.9242 - val_
accuracy: 0.8707
Epoch 25/30
750/750 [=====] - 4s 5ms/step - loss: 0.0275 - accuracy: 0.9906 - val_loss: 1.0028 - val_
accuracy: 0.8673
Epoch 26/30
750/750 [=====] - 4s 5ms/step - loss: 0.0380 - accuracy: 0.9877 - val_loss: 0.9444 - val_
accuracy: 0.8698
Epoch 27/30
750/750 [=====] - 4s 5ms/step - loss: 0.0253 - accuracy: 0.9924 - val_loss: 1.0067 - val_
accuracy: 0.8653
Epoch 28/30
750/750 [=====] - 4s 5ms/step - loss: 0.0252 - accuracy: 0.9920 - val_loss: 0.9689 - val_
accuracy: 0.8697
Epoch 29/30
750/750 [=====] - 3s 5ms/step - loss: 0.0243 - accuracy: 0.9914 - val_loss: 1.0763 - val_
accuracy: 0.8668
Epoch 30/30
750/750 [=====] - 4s 5ms/step - loss: 0.0350 - accuracy: 0.9893 - val_loss: 1.1159 - val_
accuracy: 0.8633
```



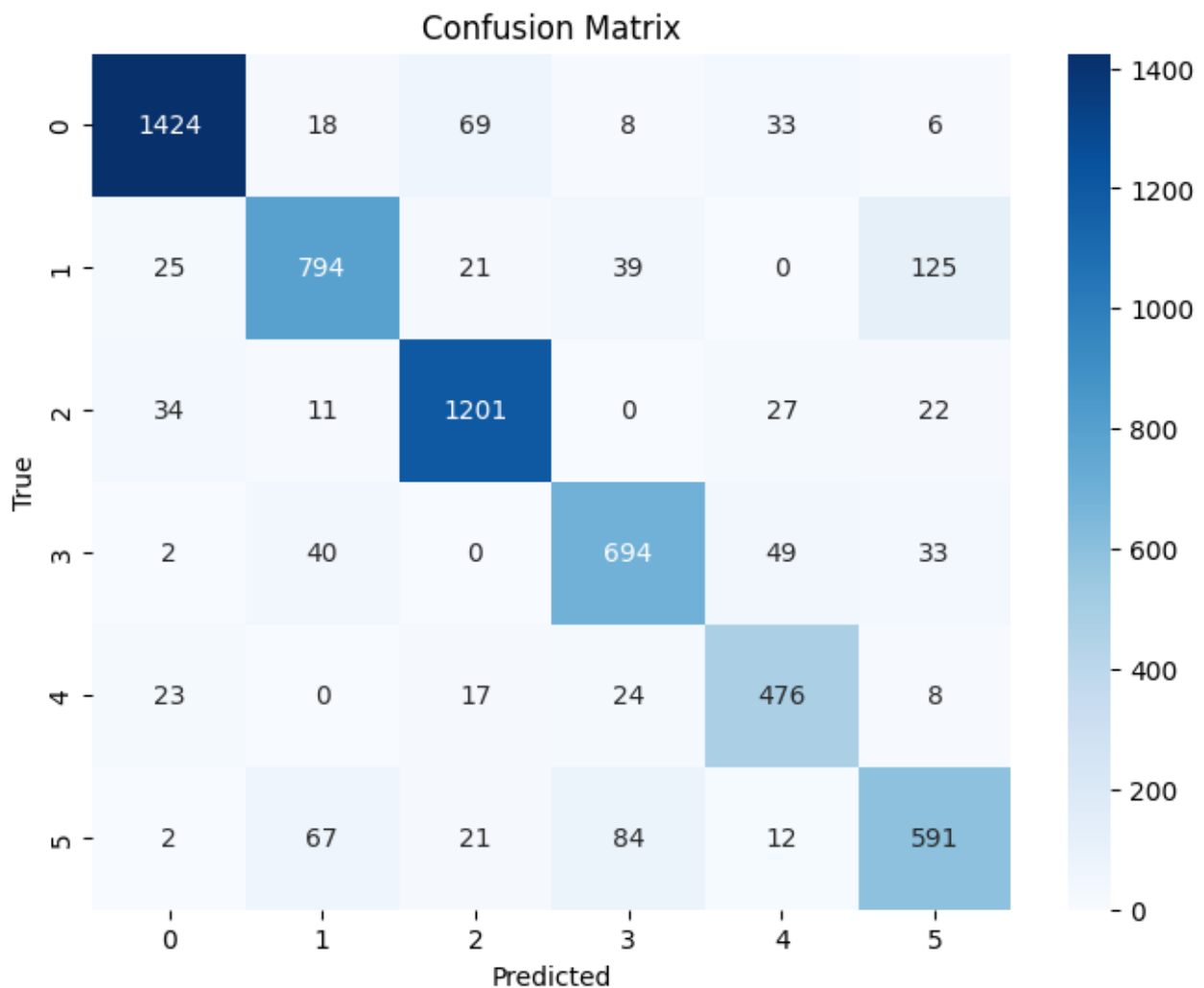
```

1 from sklearn.metrics import classification_report
2 y_pred = model.predict(X_test)
3 y_pred_labels = tf.argmax(y_pred, axis=1).numpy()
4 print(classification_report(y_test, y_pred_labels))

```

188/188 [=====] - 0s 2ms/step

	precision	recall	f1-score	support
0	0.94	0.91	0.93	1558
1	0.85	0.79	0.82	1004
2	0.90	0.93	0.92	1295
3	0.82	0.85	0.83	818
4	0.80	0.87	0.83	548
5	0.75	0.76	0.76	777
accuracy			0.86	6000
macro avg	0.84	0.85	0.85	6000
weighted avg	0.86	0.86	0.86	6000



Conclusion

We learnt the following things from our project.

1. The concept of cluster based semi-supervised learning.
2. hands on use of PCA and tsne to make cluster. (Did a lot of experimentation there)
3. Various types of clustering techniques and how they work.
4. Made a benchmarking code but it did not work well as sillouhete score is not valid for non-overlapping clusters
5. The way people point out issues on github is very ambiguous and so Encoders liek BERT, Alberta and Roberta are unable to capture that.
6. Implemented Langchain to annotation but LLMs hallucinate a lot are unable to perform well.
7. Manually labelled the clusters and got the model accuracy of **86 percent** (approx.)
Please find attached my code and other files.