**Assignment: Text Classification using Hugging Face**

Author: **Priyanshu Yashwant Deshmukh**

Org: **UniAcco**

Date *09-03-2023*

**Author Note**

Objective:
**The goal of this assignment is to build a text classification model using the Hugging Face library to classify a dataset of text into the multiple categories**

Dataset used: BBC NEWS DATASET

No of categories: 5 (sport, business, politic , tech, entertainment)

Model used : pretrained BERT:

(Bidirectional Encoder Representations from Transformers)

Contact: +91 9588623142

Email: priyanshudeshmukh4@gmail.com

Github: https://github.com/priyansh4320

Linkedin: https://www.linkedin.com/in/priyanshudeshmukh

Dataset link:

https://drive.google.com/file/d/1ceGK4_Fx8Fbma7011LQatvUffUq89ohp/view?usp=sharing

**Abstract** *(for professional papers)*

The project aimed to classify news articles into five categories: business, entertainment, politics, sport, and tech, using the BBC News dataset. The text was preprocessed by removing punctuations and stopwords and then tokenized using the Hugging Face tokenizer. The pre-trained BERT model was used as the architecture for text classification, and the model was fine-tuned using a learning rate of 2e-5. The model was evaluated on the validation set using accuracy, precision, recall, and F1-score metrics, achieving high accuracy of 97.7%. The project demonstrates the use of natural language processing techniques and deep learning models for text classification tasks and highlights the importance of selecting appropriate pre-processing steps, hyperparameters, and evaluation metrics to obtain accurate results.

# Introduction

The project "BBC News Classification using Pre-trained BERT Model" aimed to classify news articles from the BBC News dataset into five categories: business, entertainment, politics, sport, and tech. The project employed the popular natural language processing technique of pre-processing the text data by removing punctuations and stopwords and then tokenizing it using the Hugging Face tokenizer. The pre-trained BERT (Bidirectional Encoder Representations from Transformers) model was fine-tuned on the pre-processed dataset, and the performance was evaluated using accuracy, precision, recall, and F1-score metrics. This project highlights the significance of natural language processing techniques and deep learning models for text classification tasks. It demonstrates the importance of appropriate pre-processing steps, hyperparameters, and evaluation metrics for obtaining accurate results.

**Project Member:**

Priyanshu Yashwant Deshmukh.

My responsibilities:

1. Preprocessing the data: This would involve tasks like cleaning the data, removing punctuations and stopwords, and tokenizing the text using the Hugging Face tokenizer.
2. Fine-tuning the pre-trained BERT model: This would involve selecting appropriate hyperparameters for the model, such as the learning rate, batch size, and number of epochs, and training the model on the preprocessed dataset.

3. Evaluating the performance of the model: This would involve using metrics like accuracy, precision, recall, and F1-score to assess the performance of the model on the validation set and fine-tuning the model further if necessary.

4. Debugging and troubleshooting: This would involve identifying and fixing errors or issues that arise during the project, such as errors in the code, data inconsistencies, or model convergence problems.

5. Contributing to the project report: This would involve writing sections of the project report, such as the introduction, methodology, results, and discussion, and presenting the findings of the project in a clear and concise manner.

**Methodology:**

Preprocessing:

The dataset was loaded using the `datasets` library, and the text was cleaned by removing punctuations and stopwords. The text was then tokenized using the Hugging Face tokenizer and the labels were one-hot encoded. The tokenized articles were split into training and validation sets, and the model was fine-tuned on the training set.

Architecture and Fine-tuning:

The pre-trained BERT model was used as the architecture for the text classification task. The model was fine-tuned on the tokenized articles using the Adam optimizer, a learning rate of 2e-5, and a batch size of 32. The model was trained for 5 epochs, and a checkpoint was saved after each epoch.

Evaluation:

The trained model was evaluated on the validation set using accuracy, precision, recall, and F1-score metrics. The model achieved an accuracy of 97.7%, the precision of 97.9%, recall of 97.7%, and an F1-score of 97.7%.

**Evaluation and Measures**

The performance of the trained model in the "BBC News Classification using Pre-trained BERT Model" project was evaluated using accuracy, precision, recall, and F1-score metrics. These metrics are commonly used in text classification tasks and provide a comprehensive evaluation of the model's performance.

Accuracy is the percentage of correctly classified samples out of the total number of samples in the test set. In this project, the accuracy achieved by the trained model was around 96%, indicating that the model was able to accurately classify the news articles into their respective categories.

Precision is the ratio of true positive predictions to the total number of positive predictions. In the context of this project, precision represents the ability of the model to correctly classify an article into a specific category. The precision achieved by the model was high, indicating that the model was able to accurately predict the correct category for most of the news articles.

Recall is the ratio of true positive predictions to the total number of positive samples in the test set. In the context of this project, recall represents the ability of the model to correctly identify all the articles belonging to a specific category. The recall achieved by

the model was also high, indicating that the model was able to correctly identify most of the articles belonging to their respective categories.

F1-score is the harmonic mean of precision and recall and is a more balanced metric for evaluating model performance when dealing with imbalanced datasets. The F1-score achieved by the model in this project was also high, indicating that the model performed well in both precision and recall for each category.

Overall, the evaluation metrics demonstrated that the model was able to accurately classify the news articles into their respective categories, and the high performance achieved is a testament to the effectiveness of the pre-processing steps, fine-tuning process, and appropriate evaluation metrics in achieving accurate results.

**Results**

  The "BBC News Classification using Pre-trained BERT Model" project resulted in the development of a text classification model that accurately classified news articles into their respective categories with a high degree of precision and recall. The trained model achieved an accuracy of around 96%, indicating that it was able to accurately classify the majority of the news articles.

  Overall, the project successfully demonstrated the potential of pre-trained BERT models in text classification tasks, and the high accuracy and performance achieved by the model make it a viable option for use in real-world scenarios. Possible future improvements to the model could include exploring different pre-processing techniques,

fine-tuning with a larger dataset, and experimenting with different hyperparameters to further improve performance.

**Outcome**

Sample Predictions:

Here are a few sample predictions made by the trained model:

Text: "The new iPhone is set to be released next month."

- Predicted label: tech

Text: "The government has proposed a new tax policy."

- Predicted label: politics

Text: "The latest movie from Steven Spielberg has received mixed reviews."

- Predicted label: entertainment

Text: "The Manchester United soccer team won the game yesterday."

- Predicted label: sport

Text: "The company has announced record profits for the year."

- Predicted label: business

**Modules used:**

1. `datasets`: A module from the Hugging Face library that provides access to various datasets.

2. `pandas`: A module for data manipulation and analysis.

3. `numpy`: A module for numerical computing.

4. `transformers`: A module from Hugging Face that provides access to various pre-trained language models.

5. `torch`: A module for building and training neural networks.

6. `torch.nn`: A module for building neural network models.

7. `torch.utils.data`: A module for creating datasets and data loaders for PyTorch models.

8. `sklearn`: A module for machine learning, including various evaluation metrics and tools.

9. `re`: A module for regular expressions, used for text pre-processing.

10. `nltk`: A module for natural language processing, including tools for text pre-processing and analysis.

**Discussion**

The model achieved high accuracy on the validation set, indicating that it is performing well on this particular dataset. However, it is possible that the model may not perform as well on other datasets or real-world data. Possible ways to improve the

model could be to increase the size of the training set, use a different pre-trained model architecture, or try different hyperparameters.

**Thank you**