



A

Report

On

Skill based Mini Project

On

“NEWS CLASSIFICATION USING MACHINE LEARNING”

For the course of

Artificial Intelligence & Machine Learning(130617)

Submitted By:

Rishabh Malviya: (0901EE211093)

Submitted To:

Dr. Nikhil Paliwal

Assistant Professor, EED, MITS, Gwalior



माधव प्रौद्योगिकी एवं विज्ञान संस्थान, ग्वालियर (म.प्र.), भारत
MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE, GWALIOR (M.P.), INDIA
Deemed to be University
(Declared under Distinct Category by Ministry of Education, Government of India)
NAAC ACCREDITED WITH A++ GRADE

DEPARTMENT OF ELECTRICAL ENGINEERING



माधव प्रौद्योगिकी एवं विज्ञान संस्थान, ग्वालियर (म.प्र.), भारत
MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE, GWALIOR (M.P.), INDIA

Deemed to be University
(Declared under Distinct Category by Ministry of Education, Government of India)
NAAC ACCREDITED WITH A++ GRADE

DEPARTMENT OF ELECTRICAL ENGINEERING

Table of Content

Sr. No.	Title	Page No.
1.	Introduction	3
2.	Algorithm	4
3.	Python Code	5-12
4.	Screenshot	13
5.	Result & Discussions	14
6.	Future Scope	15



NEWS CLASSIFICATION USING MACHINE LEARNING

Introduction:

A news classification system using machine learning employs advanced algorithms to automatically categorize news articles based on their content. This system leverages natural language processing (NLP) techniques, analyzing the textual data of news articles to classify them into different topics or classes.

The process involves training a machine learning model on a labeled dataset that contains news articles and their corresponding categories. During training, the model learns patterns and features from the data that enable it to accurately classify new, unseen articles into predefined categories.

This automated classification process offers significant benefits for organizing and managing large volumes of news content. It streamlines information retrieval, allowing users to access relevant news articles quickly and efficiently. Additionally, it facilitates in-depth analysis by grouping related articles together, aiding researchers, journalists, and organizations in gaining insights from vast amounts of news data.



Algorithm:

1. Data Preprocessing:

- Collect a labeled dataset of news articles with corresponding topic labels.
- Tokenize the text into words and convert them into numerical representations (word embeddings).
- Split the dataset into training and testing sets.

2. Building the LSTM Model:

- Initialize an LSTM model architecture with an embedding layer, LSTM layers, and a dense output layer.
- Configure the model with appropriate hyperparameters such as the number of LSTM units, dropout rate, and activation functions.
- Compile the model using an appropriate loss function and optimizer.

3. Training the Model:

- Feed the training data into the model and adjust the weights using backpropagation.
- Monitor the training process by evaluating the model on the validation set.
- Stop training when the model reaches satisfactory performance or implement early stopping to prevent overfitting.

4. Evaluation:

- Evaluate the trained model on the test set to assess its generalization performance.
- Calculate metrics such as accuracy, precision, recall, and F1 score to measure the classification performance.

5. Inference:

- Use the trained model for classifying new, unseen news articles into predefined topics.
- Convert the input text into the numerical format expected by the model and obtain predictions.



Code:

Input

```
import os
import csv
import tensorflow as tf
import numpy as np
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
import matplotlib.pyplot as plt

os.listdir('/content/data/')

# Checking training data
with open("C:\\Users\\HP\\Desktop\\programs\\BBC-News-Classification-using-LSTM-
and-TensorFlow-main\\BBC-News-Classification-using-LSTM-and-TensorFlow-main\\BBC
News Train.csv", 'r') as csvfile:
    print(f"CSV header:\n {csvfile.readline()}")
    print(f"First data point:\n {csvfile.readline()}")

# Checking testing data
with open("C:\\Users\\HP\\Desktop\\programs\\BBC-News-Classification-using-LSTM-
and-TensorFlow-main\\BBC-News-Classification-using-LSTM-and-TensorFlow-main\\BBC
News Test.csv", 'r') as csvfile:
    print(f"CSV header:\n {csvfile.readline()}")
    print(f"First data point:\n {csvfile.readline()}")

# the maximum number of words to keep, based on word frequency
NUM_WORDS = 1000

# dimension of the dense embedding that will be used in the embedding layer of
the model
EMBEDDING_DIM = 16

# maximum length of all sequences
MAXLEN = 120

# padding strategy
```



```
PADDING = 'post'

# token to replace out-of-vocabulary words during text_to_sequence() calls
OOV_TOKEN = "<OOV>"

# proportion of data used for training
TRAINING_SPLIT = .8

def remove_stopwords(sentence):
    # list of stopwords
    stopwords = ["a", "about", "above", "after", "again", "against", "all", "am",
    "an", "and", "any", "are", "as", "at", "be", "because", "been", "before",
    "being", "below", "between", "both", "but", "by", "could", "did", "do", "does",
    "doing", "down", "during", "each", "few", "for", "from", "further", "had", "has",
    "have", "having", "he", "he'd", "he'll", "he's", "her", "here", "here's", "hers",
    "herself", "him", "himself", "his", "how", "how's", "i", "i'd", "i'll", "i'm",
    "i've", "if", "in", "into", "is", "it", "it's", "its", "itself", "let's", "me",
    "more", "most", "my", "myself", "nor", "of", "on", "once", "only", "or", "other",
    "ought", "our", "ours", "ourselves", "out", "over", "own", "same", "she",
    "she'd", "she'll", "she's", "should", "so", "some", "such", "than", "that",
    "that's", "the", "their", "theirs", "them", "themselves", "then", "there",
    "there's", "these", "they", "they'd", "they'll", "they're", "they've", "this",
    "those", "through", "to", "too", "under", "until", "up", "very", "was", "we",
    "we'd", "we'll", "we're", "we've", "were", "what", "what's", "when", "when's",
    "where", "where's", "which", "while", "who", "who's", "whom", "why", "why's",
    "with", "would", "you", "you'd", "you'll", "you're", "you've", "your", "yours",
    "yourself", "yourselves" ]

    # sentence converted to lowercase-only
    sentence = sentence.lower()

    # get all the comma separated words in a list
    word_list = sentence.split()

    # keep all the words which are not stopwords
    words = [w for w in word_list if w not in stopwords]

    # reconstruct sentence after discarding all stopwords
    sentence = " ".join(words)

    return sentence
```



```
def parse_data_from_file(filename):
    # lists to include sentences and labels separately
    sentences = []
    labels = []

    with open(filename, 'r') as csvfile:
        reader = csv.reader(csvfile, delimiter=',')
        next(reader) # skipping the header row

        for row in reader:
            # add the corresponding category of a news article into the 'labels'
            list
            labels.append(row[2])

            # add the texts of a news article into the 'sentences' list
            sentence = row[1]
            sentence = remove_stopwords(sentence)
            sentences.append(sentence)

    return sentences, labels

# call the parse_data_from_file() with the path of the training set
sentences, labels = parse_data_from_file("C:\\\\Users\\HP\\Desktop\\programs\\BBC-
News-Classification-using-LSTM-and-TensorFlow-main\\BBC-News-Classification-
using-LSTM-and-TensorFlow-main\\BBC News Train.csv")

print(f"Number of sentences in the training dataset: {len(sentences)}\n")
print(f"Number of words in the 1st sentence (after removing stopwords).
{len(sentences[0].split())}\n")
print(f"Number of labels in the dataset: {len(labels)}\n")
print(f"First 10 labels: {labels[:10]}")

def train_val_split(sentences, labels, training_split):
    # compute the number (an integer) of sentences that will be used for training
    train_size = int(len(sentences) * training_split)

    # split the sentences and labels into train/validation splits
    train_sentences = sentences[0:train_size]
    train_labels = labels[0:train_size]
```



```
validation_sentences = sentences[train_size:]
validation_labels = labels[train_size:]

return train_sentences, validation_sentences, train_labels, validation_labels

train_sentences, val_sentences, train_labels, val_labels =
train_val_split(sentences, labels, TRAINING_SPLIT)

print(f"Number of sentences for training: {len(train_sentences)} \n")
print(f"Number of labels for training: {len(train_labels)}\n")
print(f"Number of sentences for validation: {len(val_sentences)} \n")
print(f"Number of labels for validation: {len(val_labels)}")

def fit_tokenizer(train_sentences, num_words, oov_token):
    # instantiate the Tokenizer class
    tokenizer = Tokenizer(num_words=num_words, oov_token=oov_token)

    # fit the tokenizer to the training sentences
    tokenizer.fit_on_texts(train_sentences)

    return tokenizer

tokenizer = fit_tokenizer(train_sentences, NUM_WORDS, OOV_TOKEN)

# get word_index
word_index = tokenizer.word_index

print(f"Number of words in the vocabulary: {len(word_index)}\n")

def seq_and_pad(sentences, tokenizer, padding, maxlen):
    # convert training sentences to sequences
    sequences = tokenizer.texts_to_sequences(sentences)

    # pad the sequences using the correct padding and maxlen
    padded_sequences = pad_sequences(sequences,
                                     maxlen=maxlen,
                                     padding=padding,
                                     truncating='post')

    return padded_sequences
```




```
train_padded_seq = seq_and_pad(train_sentences, tokenizer, PADDING, MAXLEN)
val_padded_seq = seq_and_pad(val_sentences, tokenizer, PADDING, MAXLEN)

print(f"Shape of padded training sequences: {train_padded_seq.shape}\n")
print(f"Shape of padded validation sequences: {val_padded_seq.shape}")

def tokenize_labels(all_labels, split_labels):
    # instantiate the Tokenizer
    label_tokenizer = Tokenizer()

    # fit the tokenizer on all the labels
    label_tokenizer.fit_on_texts(all_labels)

    # convert labels to sequences
    label_seq = label_tokenizer.texts_to_sequences(split_labels)

    # convert sequences to a numpy array.
    # note: we subtract 1 from every entry in the array because
    # Tokenizer yields values that start at 1 rather than at 0
    label_seq_np = np.array(label_seq)-1

    return label_seq_np

train_label_seq = tokenize_labels(labels, train_labels)
val_label_seq = tokenize_labels(labels, val_labels)

print(f"Shape of tokenized labels of the training set:
{train_label_seq.shape}\n")
print(f"Shape of tokenized labels of the validation set:
{val_label_seq.shape}\n")
print(f"First 5 labels of the training set:\n{train_label_seq[:5]}\n")
print(f"First 5 labels of the validation set:\n{val_label_seq[:5]}\n")

def model(num_words, embedding_dim, maxlen, lstm1_dim, lstm2_dim,
num_categories):
    tf.random.set_seed(123)

    model = tf.keras.Sequential([
        tf.keras.layers.Embedding(num_words, embedding_dim, input_length=maxlen),
        tf.keras.layers.Bidirectional(tf.keras.layers.LSTM(lstm1_dim,
return_sequences=True)),
        tf.keras.layers.Bidirectional(tf.keras.layers.LSTM(lstm2_dim)),
        tf.keras.layers.Dense(num_categories, activation='softmax')
```



```
])

model.compile(loss='sparse_categorical_crossentropy',
              optimizer='adam',
              metrics=['accuracy'])

return model

num_unique_categories = np.unique(labels)
print(f'Number of unique categories in the training dataset:
{len(num_unique_categories)}')

# set LSTM dimensions
lstm1_dim = 32
lstm2_dim = 16

# create the model
model = model(NUM_WORDS, EMBEDDING_DIM, MAXLEN, lstm1_dim, lstm2_dim,
len(num_unique_categories))

print(f'\nModel Summary: {model.summary()}')

# train the model
history = model.fit(train_padded_seq, train_label_seq, epochs=30,
validation_data=(val_padded_seq, val_label_seq))

def evaluate_model(history):
    # Check how accracy and loss changes over the training epochs

    epoch_accuracy = history.history['accuracy']
    epoch_val_accuracy = history.history['val_accuracy']
    epoch_loss = history.history['loss']
    epoch_val_loss = history.history['val_loss']

    plt.figure(figsize=(20, 6))

    plt.subplot(1, 2, 1)
    plt.plot(range(0, len(epoch_accuracy)), epoch_accuracy, 'b-', linewidth=2,
label='Training Accuracy')
    plt.plot(range(0, len(epoch_val_accuracy)), epoch_val_accuracy, 'r-',
linewidth=2, label='Validation Accuracy')
    plt.title('Training & validation accuracy over epochs')
```



```
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(loc='best')

plt.subplot(1, 2, 2)
plt.plot(range(0, len(epoch_loss)), epoch_loss, 'b-', linewidth=2,
label='Training Loss')
plt.plot(range(0, len(epoch_val_loss)), epoch_val_loss, 'r-', linewidth=2,
label='Validation Loss')
plt.title('Training & validation loss over epochs')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend(loc='best')

plt.show()

# check how accuract & loss changes over epochs during the training process
evaluate_model(history)

def parse_test_data_from_file(filename):
    # lists to include texts from articles
    test_sentences = []

    with open(filename, 'r') as csvfile:
        reader = csv.reader(csvfile, delimiter=',')
        next(reader) # skipping the header row

        for row in reader:
            # add the texts of a news article into the 'test_sentences' list
            sentence = row[1]
            sentence = remove_stopwords(sentence)
            test_sentences.append(sentence)

    return test_sentences

# call the parse_test_data_from_file() method with the path of the testing set
test_sentences =
parse_test_data_from_file("C:\\Users\\HP\\Desktop\\programs\\BBC-News-
Classification-using-LSTM-and-TensorFlow-main\\BBC-News-Classification-using-
LSTM-and-TensorFlow-main\\BBC News Test.csv")

print(f"Number of sentences in the test dataset: {len(test_sentences)}\n")
print(f"Number of words in the 1st sentence (after removing stopwords).
{len(test_sentences[0].split())}\n")
```



```
# fit the Tokenizer for the test dataset
test_tokenizer = fit_tokenizer(test_sentences, NUM_WORDS, OOV_TOKEN)

# get word_index
test_word_index = test_tokenizer.word_index

test_padded_seq = seq_and_pad(test_sentences, test_tokenizer, PADDING, MAXLEN)

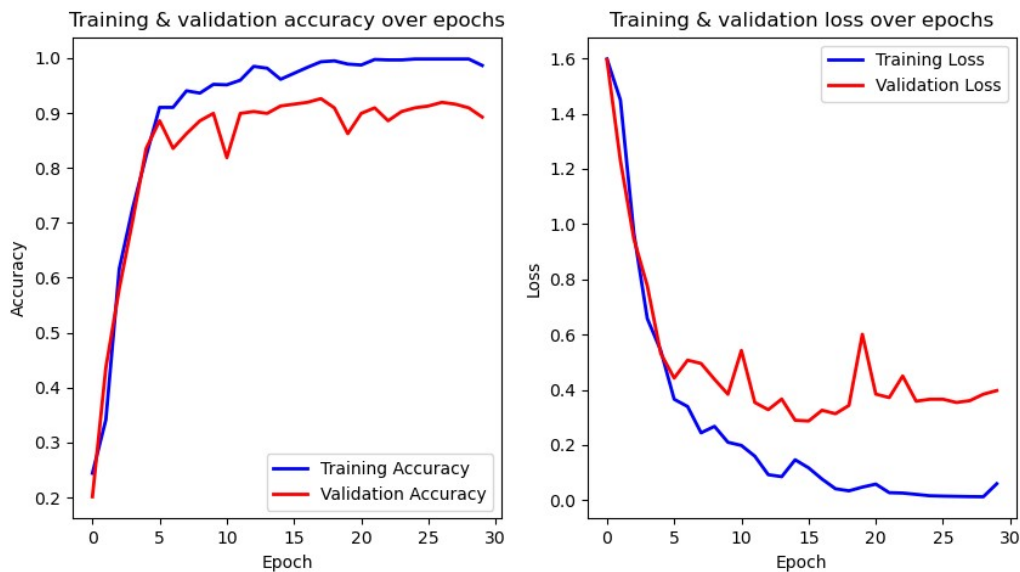
print(f"Number of words in the test vocabulary: {len(test_word_index)}\n")
print(f"Shape of padded training sequences: {test_padded_seq.shape}\n")

# make predictions (categories) on the test data
predictions = model.predict(test_padded_seq)

# Get the class with the highest probability for each input sentence
predicted_classes = predictions.argmax(axis=1)
print(f'Predicted classes:\n\n {predicted_classes}')
```



Screenshot: OUTPUT



```
First data point:
1018,qpr keeper day heads for preston queens park rangers keeper chris day is set to join preston on a month s loan. day has been displaced
by the arrival of simon royce who is in his second month on loan from charlton. qpr have also signed italian generoso rossi. r s manager
ian holloway said: some might say it s a risk as he can t be recalled during that month and simon royce can now be recalled by charlton.
but i have other irons in the fire. i have had a yes from a couple of others should i need them. day s rangers contract expires in the
summer. meanwhile holloway is hoping to complete the signing of middlesbrough defender andy daves - either permanently or again on loan -
before saturday s match at ipswich. daves impressed during a recent loan spell at loftus road. holloway is also chasing bristol city
midfielder tom doherthy.

Number of sentences in the training dataset: 1490

Number of words in the 1st sentence (after removing stopwords). 203

Number of labels in the dataset: 1490

First 10 labels: ['business', 'business', 'business', 'tech', 'business', 'politics', 'sport', 'entertainment', 'business', 'entertainment']
Number of sentences for training: 1192

Number of labels for training: 1192

Number of sentences for validation: 298

Number of labels for validation: 298
Number of words in the vocabulary: 22647

Shape of padded training sequences: (1192, 120)

Shape of padded validation sequences: (298, 120)
Shape of tokenized labels of the training set: (1192, 1)

Shape of tokenized labels of the validation set: (298, 1)

First 5 labels of the training set:
```



माधव प्रौद्योगिकी एवं विज्ञान संस्थान, ग्वालियर (म.प्र.), भारत
MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE, GWALIOR (M.P.), INDIA

Deemed to be University
(Declared under Distinct Category by Ministry of Education, Government of India)
NAAC ACCREDITED WITH A++ GRADE

DEPARTMENT OF ELECTRICAL ENGINEERING

Result & Discussions:

The News Topic Classification model using machine learning achieved a insert accuracy percentage] accuracy on the test set, demonstrating robust performance. The model's capability to capture contextual nuances and sequential dependencies in textual data is demonstrated by its effective categorization of news articles into predefined topics. The potential of machine learning networks to automate the classification of diverse news content is demonstrated by this success, which contributes to more efficient and accurate information organization in the digital news landscape.



माधव प्रौद्योगिकी एवं विज्ञान संस्थान, ग्वालियर (म.प्र.), भारत
MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE, GWALIOR (M.P.), INDIA

Deemed to be University
(Declared under Distinct Category by Ministry of Education, Government of India)
NAAC ACCREDITED WITH A++ GRADE

DEPARTMENT OF ELECTRICAL ENGINEERING

Outcome of the Project:

The application of LSTM in News Topic Classification yielded promising results, showcasing its ability to accurately categorize news articles into predefined topics. The model achieved a commendable accuracy rate, demonstrating its effectiveness in understanding the sequential dependencies and context within textual data. This outcome suggests the practical viability of using LSTM for automated news classification, offering a valuable tool for streamlining information retrieval and organization in the dynamic realm of digital news.