	https://www.kaggle.com/code/kredy10/simple-lstm-for-text-classification/data
	2. Import the required libraries import os import re import pandas as pd
	<pre>import numpy as np import nltk from nltk.corpus import stopwords from nltk.stem import WordNetLemmatizer from wordcloud import WordCloud import matplotlib.pyplot as plt import tensorflow as tf from tensorflow.keras.models import Sequential from tensorflow.keras.layers import Dense, LSTM, Dropout, Embedding from tensorflow.keras.callbacks import EarlyStopping from tensorflow.keras.preprocessing.text import Tokenizer import keras from sklearn.preprocessing import LabelEncoder</pre>
n []:	<pre>from sklearn.feature_extraction.text import TfidfVectorizer from sklearn.model_selection import train_test_split from google.colab import drive #Mount and access drive drive.mount('/content/drive', force_remount=True) os.chdir('/content/drive/My Drive')</pre>
	mounted at /content/drive Change successful. 3.Read the dataset and do pre-processing
n []: ut[]:	<pre>spam_df = pd.read_csv(filepath_or_buffer='Dataset-3_Spam.csv', delimiter=',',encoding='latin-1') spam_df.head() v1</pre>
	1 ham Ok lar Joking wif u oni NaN NaN NaN 2 spam Free entry in 2 a wkly comp to win FA Cup fina NaN NaN 3 ham U dun say so early hor U c already then say NaN NaN NaN 4 ham Nah I don't think he goes to usf, he lives aro NaN NaN NaN
n []: ut[]: n []:	<pre>#List the column names spam_df.columns Index(['v1', 'v2', 'Unnamed: 3', 'Unnamed: 4'], dtype='object') #Drop the unnamed columns</pre>
ut[]: n []:	<pre>spam_df.drop(columns=['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], axis=1, inplace=True) spam_df.columns Index(['v1', 'v2'], dtype='object') #Print_the number of rows in the dataset</pre>
ut[]: n []:	<pre>spam_df.shape (5572, 2) #Get the summary statistics of the dataset spam_df.describe()</pre>
ut[]:	v1 v2 count 5572 5572 unique 2 5169 top ham Sorry, I'll call later freq 4825 30
n []: ut[]:	<pre>#Check for null values spam_df.isna().sum()</pre>
n []: ut[]:	<pre>dtype: int64 #Check for duplicated rows spam_df.duplicated().sum()</pre>
n []: ut[]:	<pre>#Remove the duplicated rows spam_df = spam_df.drop_duplicates() spam_df.duplicated().sum()</pre>
n []:	<pre>#Display the count of spam and ham labels #Stratified-split is required spam_df['v1'].hist(bins=3) <matplotlib.axessubplots.axessubplot 0x7f92aa5f2f90="" at=""> 4000 2000 1000 1000 1000 1000 1000 100</matplotlib.axessubplots.axessubplot></pre>
[]:	<pre>def wordcloud_vis(column): mostcommon = nltk.FreqDist(spam_df[column]).most_common(100) wordcloud = WordCloud(width=1600, height=800, background_color='white').generate(str(mostcommon)) fig = plt.figure(figsize=(30,10), facecolor='white') plt.imshow(wordcloud) #, interpolation="bilinear") plt.axis('off') plt.show()</pre>
ı []:	#Plot the word-cloud before removing stopwords, performing lemmatization wordcloud_vis('v2') Cine
n []:	msg spell miss awarded a speak always on the speak always only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space spain with thanks only the letters and spaces spain did aluna text is a space space in the spac
t[]:	spam_df['alpha_text'] = spam_df['v2'].apply(lambda x: re.sub(r'[^a-zA-Z]+', '', x.lower())) v1
]:	<pre>#Remove stop-words nltk.download('stopwords') spam_df['imp_text'] = spam_df['alpha_text'].apply(lambda x : ' '.join([word for word in x.split() if not word in set(stopwords.words('english') spam_df.head() [nltk_data] Downloading package stopwords to /root/nltk_data [nltk_data] Unringing corpore(stopwords ring)</pre>
t[]:	Unzipping corpora/stopwords.zip. v1 v2 alpha_text imp_text o ham Go until jurong point, crazy Available only go until jurong point crazy available only in go jurong point crazy available bugis n great ham Ok lar Joking wif u oni ok lar joking wif u oni ok lar joking wif u oni spam Free entry in 2 a wkly comp to win FA Cup fina free entry in a wkly comp to win fa cup final free entry wkly comp win fa cup final tkts st
ı []:	ham V dun say so early hor U c already then say u dun say so early hor u c already then say u dun say early hor u c already say ham Nah I don't think he goes to usf, he lives aro nah i dont think he goes to usf he lives aroun nah dont think goes usf lives around though #Tokenize the data def tokenize(data):
ıt[]:	<pre>generated_token = list(data.split()) return generated_token spam_df['token_text'] = spam_df['imp_text'].apply(lambda x: tokenize(x)) spam_df.head() v1</pre>
	1 ham Ok lar Joking wif u oni [ok, lar, joking, wif, u, oni] 2 spam Free entry in 2 a wkly comp to win FA Cup fina free entry in a wkly comp to win fa cup final 3 ham U dun say so early hor U c already then say u dun say so early hor u c already then say 4 ham Nah I don't think he goes to usf, he lives aro nah i dont think he goes to usf he lives aroun Ok lar joking wif u oni free entry wkly comp win fa cup final tkts st [free, entry, wkly, comp, win, fa, cup, final, [free, entry, wkly, comp, win, fa, cup, final, u dun say early hor u c already say u dun say early hor u c already say nah dont think goes usf lives around though [nah, dont, think, goes, usf, lives, around, t
	<pre>#Perform lemmatization nltk.download('wordnet')</pre>
[]:	<pre>nltk.download('omw-1.4') lemmatizer = WordNetLemmatizer() def lemmatization(list_of_words): lemmatized_list = [lemmatizer.lemmatize(word) for word in list_of_words] return lemmatized_list</pre>
	<pre>lemmatizer = WordNetLemmatizer() def lemmatization(list_of_words): lemmatized_list = [lemmatizer.lemmatize(word) for word in list_of_words] return lemmatized_list spam_df['lemmatized_text'] = spam_df['token_text'].apply(lambda x: lemmatization(x)) spam_df.head() [nltk_data] Downloading package wordnet to /root/nltk_data [nltk_data] Package wordnet is already up-to-date! [nltk_data] Downloading package omw-1.4 to /root/nltk_data</pre>
	lemmatizer = WordNetLemmatizer() def lemmatization(list_of_words): lemmatized_list = [lemmatizer.lemmatize(word) for word in list_of_words] return lemmatized_list spam_df['lemmatized_list spam_df['lemmatized_lemmatized_list spam_df['lemmatized_lemmatized_lemmatized_list spam_df['lemmatized_lemmatized_lemmatized_lemmatized_list spam_df['lemmatized_lem
ıt[]:	Lemmatizer wordNetLemmatizer() def lemmatization(list_of_words);
nt[]:	lemmatized_list = [lemmatizet] deference of provords] return lemmatized_list = [spam_dff'(token_text'].apply(lambda x: lemmatization(x)) spam_df'.head() [nltk_data] Downloading package wordnet to /root/nltk_data [nltk_data] Downloading package wordnet to /root/nltk_data [nltk_data] Downloading package wordnet is already up-to-date! [nlt
t[]:	Demantizer Bornthettemmatizer Order Demantizer Order Demantizer Deman
t[]:	Item
t[]:	To control to the con
t[]:	International Continues (Continues of State Continues of State Conti
t[]: t[]:	Segretaria - Morrison Consecution (1) Service 1 Morrison (1) Service 1 Morri
t[]: t[]:	Service Servic
t[]: t[]: t[]:	Control Cont
it[]: it[]: it[]: it[]:	The control of the co
t[]: t[]: t[]: t[]:	The control co
it[]:	## Provided to Company (1997) 1997
it[]:	Description Company
t[]: t[]: t[]: t[]: t[]: t[]: t[]: t[]:	Description of the property of
	### 1997 Part of the content of th
	### 1985 1985
	The second content of the content of
	The content is a property of the content is
	The state of the s