

In [3]:

```
from google.colab import drive
drive.mount('/content/drive/')
```

Drive already mounted at /content/drive/; to attempt to forcibly remount, call drive.mount("/content/drive/", force_remount=True).

In [4]:

```
import nltk
import pandas as pd
```

In [5]:

```
nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

Out[5]:

True

In [6]:

```
from nltk.tokenize import word_tokenize, sent_tokenize
filepath = '/content/drive/My Drive/English-Telugu/train.en'
corpus = open(filepath, 'r').read()
words = nltk.word_tokenize(corpus)
print("The number of tokens is", len(words))
average_tokens = round(len(words)/75000)
print("The average number of tokens per sentence is", average_tokens)
unique_tokens = set(words)
print("The number of unique tokens are", len(unique_tokens))
```

The number of tokens is 1809312
The average number of tokens per sentence is 24
The number of unique tokens are 21095

In [7]:

```
from nltk.tokenize import word_tokenize, sent_tokenize
filepath = nltk.data.find('/content/drive/My Drive/English-Telugu/train.te')
corpus = open(filepath, 'r').read()
words = nltk.word_tokenize(corpus)
print("The number of tokens is", len(words))
average_tokens = round(len(words)/75000)
print("The average number of tokens per sentence is", average_tokens)
unique_tokens = set(words)
print("The number of unique tokens are", len(unique_tokens))
```

The number of tokens is 1030924
The average number of tokens per sentence is 14
The number of unique tokens are 106016

In [8]:

```
f= open("/content/drive/My Drive/English-Telugu/train.en")
en=f.readlines()
len(en)
```

Out[8]:

75000

In [9]:

```
en[:6]
```

['we just party , and we can do whatever we want .\n',
 'and ziza the son of shiphi , the son of allon , the son of jedaiah , the son of shimri
 , the son of shemaiah;\n',
 'now a poor wise man was found in it , and he by his wisdom delivered the city; yet no m
an remembered that same poor man .\n',
 "and the child grew , and she brought him unto pharaoh's daughter , and he became her so
n . and she called his name moyses: and she said , because i drew him out of the water .\n",
 '- i was at the jeweler and nobody came .\n',
 'and the consecrated things were six hundred oxen and three thousand sheep .\n']

```
f= open("/content/drive/My Drive/English-Telugu/train.te")
te=f.readlines()
len(te)
```

75000

```
import re
def remove_punc(x):
    return re.sub('[!#?,.:;"'\n]', '', x)
```

```
for i in range(len(te)):
    te[i]=remove_punc(te[i])
```

```
te[:6]
```

[illegible]

```
for i in range(len(en)):
    en[i]=remove_punc(en[i])
```

```
en[:6]
```

['we just party and we can do whatever we want ',
'and ziza the son of shiphi the son of allon the son of jedaiah the son of shimri the son of shemaiah',
'now a poor wise man was found in it and he by his wisdom delivered the city yet no man remembered that same poor man ',
"and the child grew and she brought him unto pharaoh's daughter and he became her son and she called his name moyses and she said because i drew him out of the water ",
'- i was at the jeweler and nobody came '

'and the consecrated things were six hundred oxen and three thousand sheep ']

In [16]:

```
words_en = []
for i in en:
    for word in i.split():
        words_en.append(word)
words_en[:10]
```

Out[16]:

```
['we', 'just', 'party', 'and', 'we', 'can', 'do', 'whatever', 'we', 'want']
```

In [17]:

```
words_te = []
for i in te:
    for word in i.split():
        words_te.append(word)
words_te[:10]
```

Out[17]:

[illegible]

In [18]:

```
from collections import Counter
english_words_counts = Counter(words_en)
telugu_words_counts = Counter(words_te)
```

In [19]:

```
import operator
english_words_counts = sorted(english_words_counts.items(), key = operator.itemgetter(1),
reverse = True)
telugu_words_counts = sorted(telugu_words_counts.items(), key = operator.itemgetter(1), r
everse = True)
```

In [20]:

```
english_words_counts[:10]
```

Out[20]:

```
[('the', 113593),
 ('and', 80265),
 ('of', 62955),
 ('to', 34296),
 ('in', 24579),
 ('that', 21052),
 ('you', 20618),
 ('he', 19632),
 ('i', 19186),
 ('a', 17782)]
```

In [21]:

```
telugu words counts[:10]
```

Out[21]:

```
[('■■■■', 12085),
 ('■■■■■■■■', 9127),
 ('■■', 9006),
 ('■■', 8842),
 ('■', 8681),
 ('■■■■', 8569),
 ('■■■■', 6680),
 ('■■', 5784),
 ('■■■■', 5701),
 ('■■■■', 5691)]
```

In [22]:

```
maxlen_english = -1
for doc in en:
    tokens = nltk.word_tokenize(doc)
    if(maxlen_english < len(tokens)):
        maxlen_english = len(tokens)
print("The maximum number of words in any document = ", maxlen_english)
```

The maximum number of words in any document = 74

In [23]:

```
maxlen_telugu = -1
for doc in te:
    tokens = nltk.word_tokenize(doc)
    if(maxlen_telugu < len(tokens)):
        maxlen_telugu = len(tokens)
print("The maximum number of words in any document = ", maxlen_telugu)
```

The maximum number of words in any document = 47

In [24]:

```
def tokenize_and_pad(x, maxlen):
    # a tokenizer to tokenize the words and create sequences of tokenized words
    tokenizer = Tokenizer(char_level = False)
    tokenizer.fit_on_texts(x)
    sequences = tokenizer.texts_to_sequences(x)
    padded = pad_sequences(sequences, maxlen = maxlen, padding = 'post')
    return tokenizer, sequences, padded
```

In [25]:

```
from tensorflow.keras.preprocessing.text import one_hot, Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
x_tokenizer, x_sequences, x_padded = tokenize_and_pad(en, maxlen_english)
y_tokenizer, y_sequences, y_padded = tokenize_and_pad(te, maxlen_telugu)
```

In [26]:

```
[print("The tokenized version for document\n", en[-1:][0], "\n", x_padded[-1:][0])]
```

The tokenized version for document

how could you stand for it

```
[149 391 7 316 11 16 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0]
```

Out[26]:

[None]

In [27]:

```
print("The tokenized version for document\n", te[-1:][0], "\n ", y_padded[-1:][0])
```

The tokenized version for document

```

[ 6 296 8102 887 0 0 0 0 0 0 0 0 0 0 0 0]
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0]

```

In [28]:

```

from tensorflow.keras.preprocessing.text import one_hot, Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Flatten, TimeDistributed, RepeatVector, Embedding, Input, LSTM, Conv1D, MaxPool1D, Bidirectional
from tensorflow.keras.models import Model

```

In [29]:

```

english_vocab_size = len(english_words_counts)
telugu_vocab_size = len(telugu_words_counts)

```

In [30]:

```

# Sequential Model
model = Sequential()
# embedding layer
model.add(Embedding(english_vocab_size, 256, input_length = maxlen_english, mask_zero = True))
# encoder
model.add(LSTM(256, return_sequences= True))
model.add(LSTM(128))
# decoder
# repeatvector repeats the input for the desired number of times to change
# 2D-array to 3D array. For example: (1,256) to (1,23,256)
model.add(RepeatVector(maxlen_telugu))
model.add(LSTM(256, return_sequences= True ))
model.add(LSTM(128, return_sequences= True ))
model.add(TimeDistributed(Dense(telugu_vocab_size, activation = 'softmax'))))
model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])
model.summary()

```

Model: "sequential"

| Layer (type) | Output Shape | Param # |
|------------------------------|--------------------|----------|
| embedding (Embedding) | (None, 74, 256) | 5810688 |
| lstm (LSTM) | (None, 74, 256) | 525312 |
| lstm_1 (LSTM) | (None, 128) | 197120 |
| repeat_vector (RepeatVector) | (None, 47, 128) | 0 |
| lstm_2 (LSTM) | (None, 47, 256) | 394240 |
| lstm_3 (LSTM) | (None, 47, 128) | 197120 |
| time_distributed (TimeDistri | (None, 47, 106518) | 13740822 |
| Total params: 20,865,302 | | |
| Trainable params: 20,865,302 | | |
| Non-trainable params: 0 | | |

In [31]:

```

from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(x_padded, y_padded, test_size = 0.1)

```

In [32]:

```
In [32]:
```

```
import numpy as np
y_train = np.expand_dims(y_train, axis = 2)
y_train.shape
```

```
Out[32]:
```

```
(67500, 47, 1)
```

```
In [33]:
```

```
history=model.fit(x_train, y_train, batch_size=64, validation_split= 0.1, epochs=10)
```

```
Epoch 1/10
950/950 [=====] - 496s 509ms/step - loss: 3.9208 - accuracy: 0.7
399 - val_loss: 2.5847 - val_accuracy: 0.7406
Epoch 2/10
950/950 [=====] - 478s 504ms/step - loss: 2.4850 - accuracy: 0.7
467 - val_loss: 2.5501 - val_accuracy: 0.7430
Epoch 3/10
950/950 [=====] - 478s 503ms/step - loss: 2.4294 - accuracy: 0.7
478 - val_loss: 2.5098 - val_accuracy: 0.7432
Epoch 4/10
950/950 [=====] - 478s 503ms/step - loss: 2.3647 - accuracy: 0.7
485 - val_loss: 2.4735 - val_accuracy: 0.7437
Epoch 5/10
950/950 [=====] - 478s 504ms/step - loss: 2.3074 - accuracy: 0.7
491 - val_loss: 2.4303 - val_accuracy: 0.7441
Epoch 6/10
950/950 [=====] - 478s 503ms/step - loss: 2.2393 - accuracy: 0.7
495 - val_loss: 2.3830 - val_accuracy: 0.7446
Epoch 7/10
950/950 [=====] - 477s 503ms/step - loss: 2.1560 - accuracy: 0.7
519 - val_loss: 2.3554 - val_accuracy: 0.7447
Epoch 8/10
950/950 [=====] - 479s 504ms/step - loss: 2.0994 - accuracy: 0.7
522 - val_loss: 2.3254 - val_accuracy: 0.7458
Epoch 9/10
950/950 [=====] - 478s 503ms/step - loss: 2.0433 - accuracy: 0.7
525 - val_loss: 2.3023 - val_accuracy: 0.7461
Epoch 10/10
950/950 [=====] - 478s 503ms/step - loss: 1.9857 - accuracy: 0.7
534 - val_loss: 2.2820 - val_accuracy: 0.7460
```

```
In [34]:
```

```
y_test = np.expand_dims(y_test, axis = 2)
y_test.shape
```

```
Out[34]:
```

```
(7500, 47, 1)
```

```
In [35]:
```

```
model.evaluate(x_test, y_test, batch_size=32)
```

```
235/235 [=====] - 61s 258ms/step - loss: 2.2293 - accuracy: 0.75
22
```

```
Out[35]:
```

```
[2.229292392730713, 0.7522156238555908]
```

```
In [36]:
```

```
model.save("NMT2.h5")
```

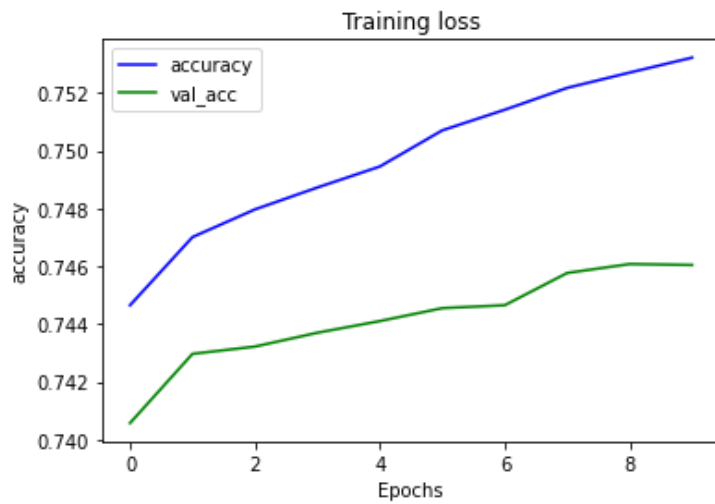
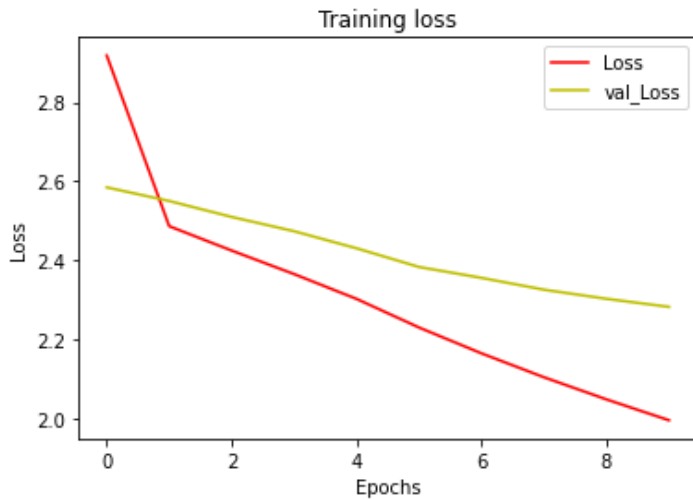
```
In [44]:
```

```
import matplotlib.pyplot as plt
loss=history.history['loss']
```

```

acc=history.history['accuracy']
val_loss=history.history['val_loss']
val_acc=history.history['val_accuracy']
epochs=range(len(loss))
plt.plot(epochs, loss, 'r')
plt.plot(epochs, val_loss, 'y')
plt.title('Training loss')
plt.xlabel("Epochs")
plt.ylabel("Loss")
plt.legend(["Loss", "val_Loss"])
plt.show()
plt.plot(epochs, acc, 'b')
plt.plot(epochs, val_acc, 'g')
plt.title('Training loss')
plt.xlabel("Epochs")
plt.ylabel("accuracy")
plt.legend(["accuracy", "val_acc"])
plt.show()

```



In [40]:

```

from google.colab import files
files.download('NMT2.h5')

```