

# Twitter Named Entity Recognition Case Study

## About

Twitter is a microblogging and social networking service on which users post and interact with messages known as "tweets". Every second, on average, around 6,000 tweets are tweeted on Twitter, corresponding to over 350,000 tweets sent per minute, 500 million tweets per day.

## Problem statement

Twitter wants to automatically tag and analyze tweets for better understanding of the trends and topics without being dependent on the hashtags that the users use. Many users do not use hashtags or sometimes use wrong or mis-spelled tags, so they want to completely remove this problem and create a system of recognizing important content of the tweets.

## Objective

Named Entity Recognition (NER) is an important subtask of information extraction that seeks to locate and recognise named entities. We need to train models that will be able to identify the various named entities.

## Data

Dataset is annotated with 10 fine-grained NER categories: person, geo-location, company, facility, product,music artist, movie, sports team, tv show and other. Dataset was extracted from tweets and is structured in CoNLL format., in English language. Containing in Text file format. The CoNLL format is a text file with one word per line with sentences separated by an empty line. The first word in a line should be the word and the last word should be the label.

```
In [1]:  
1 # imports  
2 import pandas as pd  
3 import numpy as np  
4 import matplotlib.pyplot as plt  
5 import seaborn as sns  
6 import tensorflow as tf  
7  
8 from warnings import filterwarnings  
9 filterwarnings('ignore')
```

WARNING:tensorflow:From c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\venv\lib\site-packages\keras\src\losses.py:2976: The name tf.losses.sparse\_softmax\_cross\_entropy is deprecated. Please use tf.compat.v1.losses.sparse\_softmax\_cross\_entropy instead.

```
In [2]: 1 # setting path variables
2 import os
3 root_path = os.path.abspath(os.path.join(os.getcwd(),os.pardir))
4 data_path = os.path.join(root_path,'data')
5 train_data_path = os.path.join(data_path,'wnut_16.txt.conll')
6 test_data_path = os.path.join(data_path,'wnut_16test.txt.conll')
```

## Getting the data

```
In [3]: 1 # reading the training file
2 with open(train_data_path, 'r') as f:
3     train_raw = f.read()
4 with open(test_data_path, 'r') as f:
5     test_raw = f.read()
```

```
In [4]: 1 # creating a function to format the data
2 def extract_ner_from_conll(conll_data):
3     # Split the data into sentences based on empty lines
4     sentences = [sentence.strip() for sentence in conll_data.strip().split()
5     ner_data = []
6
7     for sentence in sentences:
8         tokenised_sentence = []
9         for token_entity in sentence.split('\n'):
10             token, entity = token_entity.split('\t')
11             tokenised_sentence.append((token,entity))
12         ner_data.append(tokenised_sentence)
13
14     return ner_data
```

```
In [5]: 1 # preprocessing the raw files
2 train_data = extract_ner_from_conll(train_raw)
3 test_data = extract_ner_from_conll(test_raw)
```

```
In [6]: 1 # checking sentences after preprocessing
2 print(train_data[0])
```

```
[('@SammieLynnsMom', '0'), ('@tg10781', '0'), ('they', '0'), ('will', '0'),
('be', '0'), ('all', '0'), ('done', '0'), ('by', '0'), ('Sunday', '0'), ('true',
'st', '0'), ('me', '0'), ('*wink*', '0')]
```

# EDA

```
In [7]: 1 # number of words in the vocabulary and length of sentences in the training
2
3 sentence_lengths = list()
4 word_set = set()
5 for sentence in train_data:
6     sentence_lengths.append(len(sentence))
7     for word in sentence:
8         word_set.add(word[0])
9
10 NUM_WORDS = len(word_set)+2 # +2 to include padding and out of vocabulary
11 print(f"Number of unique words in training data (including padding and OOV) = {NUM_WORDS}")
12 print(f"Maximum sentence length = {max(sentence_lengths)}")
13 print(f"Minimum sentence length = {min(sentence_lengths)}")
```

Number of unique words in training data (including padding and OOV token) = 10588

Maximum sentence length = 39

Minimum sentence length = 1

```
In [8]: 1 # since the max sentence length is 39, we will take a length of 45 in our
2 SENTENCE_LENGTH = 45
3 # we will keep the embedding dimensions to be 50 since the number of datapoints is less than 1000
4 EMBEDDING_DIMS = 300
```

```
In [9]: 1 # number of entities
2 entity_set = set()
3 for sentence in train_data:
4     for word in sentence:
5         entity_set.add(word[1])
6
7 NUM_ENTITIES = len(entity_set)
8 print(f"Number of unique entities in training data = {NUM_ENTITIES}")
```

Number of unique entities in training data = 21

## Data preparation

```
In [10]: 1 import re
2 import string
3 punctuations = string.punctuation
```

```
In [11]: 1 # create a function to prepare the data to be fed into the model
2
3 def clean_text(text):
4     return re.sub("[^A-Za-z0-9]+", ' ', str(text).lower())
5
6 def prepare_data(text_data):
7
8     # initialize empty lists for sentences and entities
9     sentences = []
10    entities = []
11
12    for sentence in text_data:
13
14        # initialize empty lists for sentence text and corresponding entit
15        word_list = []
16        entity_list = []
17
18        for token in sentence:
19            word = token[0]
20            entity = token[1]
21            if word in punctuations:
22                continue
23            else:
24                word_list.append(clean_text(word))
25                entity_list.append(entity)
26
27        sentences.append(word_list)
28        entities.append(entity_list)
29
30    # create a single string for each sentence and entity by joining eleme
31    sentences = np.array([' '.join(sentence) for sentence in sentences])
32    entities = np.array([' '.join(entity) for entity in entities])
33
34    return (sentences,entities)
```

```
In [12]: 1 # since the datapoints in train file is very low, we will merge the dataset
2 data = []
3 for i in train_data:
4     data.append(i)
5 for i in test_data:
6     data.append(i)
7
8 xdata, ydata = prepare_data(data)
9
10 from sklearn.model_selection import train_test_split
11 xtrain, xtest, ytrain, ytest = train_test_split(xdata, ydata, test_size=0.
```

```
In [13]: 1 # checking sentences are conversion
2 print(f"Before conversion\n{' '.join([word[0] for word in data[0]])}\n")
3 print(f"After conversion\n{xdata[0]}\n")
4 print(f"Entities\n{ydata[0]}")
```

Before conversion

@SammieLynnsMom @tg10781 they will be all done by Sunday trust me \*wink\*

After conversion

sammielynnsmom tg10781 they will be all done by sunday trust me wink

Entities

0 0 0 0 0 0 0 0 0 0 0 0

## Bidirectional LSTM + CRF

```
In [14]: 1 from tensorflow.keras.layers import Input, TextVectorization, Embedding, B
2 from tensorflow.keras.models import Model
3 from tensorflow_addons.layers import CRF
4 from tensorflow_addons.losses import SigmoidFocalCrossEntropy
```

```
In [15]: 1 # adapt the vectorizer layer before modeling
2 sentence_tokenizer = TextVectorization(max_tokens=NUM_WORDS, output_sequen
3 sentence_tokenizer.adapt(xtrain)
4
5 # indexing the tokens in train and test
6 train_lstm_sentence_indexed = sentence_tokenizer(xtrain)
7 test_lstm_sentence_indexed = sentence_tokenizer(xtest)
```

WARNING:tensorflow:From c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\venv\lib\site-packages\keras\src\backend.py:873: The name tf.get\_default\_graph is deprecated. Please use tf.compat.v1.get\_default\_graph instead.

WARNING:tensorflow:From c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\venv\lib\site-packages\keras\src\utils\tf\_utils.py:492: The name tf.ragged.RaggedTensorValue is deprecated. Please use tf.compat.v1.ragged.RaggedTensorValue instead.

```
In [16]: 1 # initializing the embedding vector
2 import gensim.downloader as api
3 word2vec = api.load('word2vec-google-news-300')
4 embedding_matrix = np.zeros(shape=(NUM_WORDS, EMBEDDING_DIMS), dtype=np.float32)
5
6 for i, word in enumerate(sentence_tokenizer.get_vocabulary()):
7     try:
8         embedding_matrix[i] = word2vec[word]
9     except:
10        pass
```

```
In [17]: 1 del word2vec
```

```
In [18]: 1 # define a function to build and compile model
2 def build_lstm_model(name='bi-lstm+crf'):
3
4     # input layer for getting sentences
5     sentence_input = Input(shape=(SENTENCE_LENGTH,), dtype=tf.float32, name='sentence_input')
6
7     # creating embeddings for each token in sentence
8     embeddings = Embedding(
9         input_dim = NUM_WORDS,
10        output_dim = EMBEDDING_DIMS,
11        mask_zero = True,
12        name = 'word_embedding',
13        embeddings_initializer = tf.keras.initializers.constant(embedding_matrix))
14     (sentence_input)
15
16     # stacking two bidirectional LSTMs
17     output_sequence = Bidirectional(LSTM(50, return_sequences=True), name='output_sequence')
18     output_sequence = Bidirectional(LSTM(50, return_sequences=True), name='output_sequence')
19
20     # passing the sequence through dense Layer to compress the information
21     dense_sequence = TimeDistributed(Dense(25, activation='relu'), name='dense_sequence')
22
23     # passing the dense sequences through crf Layer
24     predicted_sequence, potentials, sequence_length, crf_kernel = CRF(NUM_TAGS).apply(
25         output_sequence)
26
27     # define the train model
28     training_model = Model(sentence_input, predicted_sequence)
29
30     # compile the model
31     training_model.compile(
32         loss=SigmoidFocalCrossEntropy(),
33         optimizer='adam'
34     )
35
36     # create an inference model
37     inference_model = Model(sentence_input, predicted_sequence)
38
39     return training_model, inference_model
```

```
In [19]: 1 # creating the traing and inferencing model
          2 lstm_training_model, lstm_inference_model = build_lstm_model()
```

WARNING:tensorflow:From c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\venv\lib\site-packages\keras\src\optimizers\\_\_init\_\_.py:309: The name tf.train.Optimizer is deprecated. Please use tf.compat.v1.train.Optimizer instead.

```
In [20]: 1 lstm_training_model.summary()
```

Model: "model"

Layer (type)	Output Shape	Param #
<hr/>		
sentence_input (InputLayer )	[None, 45]	0
word_embedding (Embedding)	(None, 45, 300)	3176400
lstm_1 (Bidirectional)	(None, 45, 100)	140400
lstm_2 (Bidirectional)	(None, 45, 100)	60400
dense (TimeDistributed)	(None, 45, 25)	2525
crf (CRF)	[None, 45), (None, 45, 21), (None,), (21, 21)]	1029
<hr/>		
Total params:	3380754 (12.90 MB)	
Trainable params:	3380754 (12.90 MB)	
Non-trainable params:	0 (0.00 Byte)	

---

```
In [21]: 1 # prepare the entity data
          2
          3 # tokenizing the entities
          4 entity_tokenizer = TextVectorization(max_tokens=NUM_ENTITIES, output_seque
          5 entity_tokenizer.adapt(ytrain)
          6 ytrain_tokenized = entity_tokenizer(ytrain)
          7 ytest_tokenized = entity_tokenizer(ytest)
          8
          9 # one hot encoding the entitiy tokens
          10 entity_ohe = Lambda(lambda x: tf.one_hot(x, NUM_ENTITIES))
          11 ytrain_tokenized = entity_ohe(ytrain_tokenized)
          12 ytest_tokenized = entity_ohe(ytest_tokenized)
```

In [22]:

```
1 # defining callbacks
2 mc = tf.keras.callbacks.ModelCheckpoint(
3     os.path.join(root_path, 'models', 'lstm.ckpt'),
4     monitor="val_loss",
5     save_best_only=True,
6     save_weights_only=True
7 )
8 es = tf.keras.callbacks.EarlyStopping(
9     monitor="val_loss",
10    min_delta=0,
11    patience=5,
12    start_from_epoch=5,
13 )
14 callbacks=[mc,es]
```

In [23]:

```
1 # fitting the Lstm+crf model
2 history_lstm = lstm_training_model.fit(
3     train_lstm_sentence_indexed, ytrain_tokenized,
4     validation_data=(test_lstm_sentence_indexed, ytest_tokenized),
5     epochs=50,
6     batch_size=16,
7     callbacks=callbacks
8 )
```

Epoch 1/50

WARNING:tensorflow:Gradients do not exist for variables ['chain\_kernel:0'] when minimizing the loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?

WARNING:tensorflow:Gradients do not exist for variables ['chain\_kernel:0'] when minimizing the loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?

WARNING:tensorflow:Gradients do not exist for variables ['chain\_kernel:0'] when minimizing the loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?

WARNING:tensorflow:Gradients do not exist for variables ['chain\_kernel:0'] when minimizing the loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?

313/313 [=====] - 36s 73ms/step - loss: 0.3820 - val\_loss: 0.1575

Epoch 2/50

313/313 [=====] - 18s 57ms/step - loss: 0.0876 - val\_loss: 0.0535

Epoch 3/50

313/313 [=====] - 17s 54ms/step - loss: 0.0367 - val\_loss: 0.0436

Epoch 4/50

313/313 [=====] - 17s 55ms/step - loss: 0.0269 - val\_loss: 0.0449

Epoch 5/50

313/313 [=====] - 19s 59ms/step - loss: 0.0235 - val\_loss: 0.0645

Epoch 6/50

313/313 [=====] - 17s 55ms/step - loss: 0.0216 - val\_loss: 0.0738

Epoch 7/50

313/313 [=====] - 17s 56ms/step - loss: 0.0206 - val\_loss: 0.1044

Epoch 8/50

313/313 [=====] - 17s 55ms/step - loss: 0.0188 - val\_loss: 0.1294

Epoch 9/50

313/313 [=====] - 18s 57ms/step - loss: 0.0212 - val\_loss: 0.1290

Epoch 10/50

313/313 [=====] - 18s 56ms/step - loss: 0.0188 - val\_loss: 0.1456

Epoch 11/50

313/313 [=====] - 18s 56ms/step - loss: 0.0180 - val\_loss: 0.1766

```
In [24]: 1 # Loading the best model
          2 lstm_training_model.load_weights(os.path.join(root_path,'models','lstm.ckpt'))
```

```
Out[24]: <tensorflow.python.checkpoint.checkpoint.CheckpointLoadStatus at 0x26c55d15e7
0>
```

```
In [25]: 1 # evaluating on validation data
          2 lstm_training_model.evaluate(test_lstm_sentence_indexed,ytest_tokenized)
```

```
40/40 [=====] - 1s 12ms/step - loss: 0.0436
```

```
Out[25]: 0.04360351338982582
```

## Inferencing the Bi-LSTM+CRF model

```
In [32]: 1 # function to infer prediction from a text in xtest randomly
          2 def infer_lstm():
          3     idx = np.random.choice(range(len(xtest)),1,replace=False)
          4     text = xtest[idx]
          5     text_tokenized = sentence_tokenizer(text)
          6     pred_labels = lstm_inference_model.predict(text_tokenized)
          7     act_labels = ytest[idx]
          8     text_len = min(len(text[0].split()),SENTENCE_LENGTH)
          9
         10    pred_labels = np.asarray(entity_tokenizer.get_vocabulary())[list(pred_
         11    result = pd.DataFrame(
         12        {
         13            'text':text[0].split(),
         14            'actual_labels':act_labels[0].split(),
         15            'predicted_labels':pred_labels[:text_len]
         16        }
         17    )
         18    display(result)
```

```
In [64]: 1 # infer random text from the validation data
          2 infer_lstm()
```

1/1 [=====] - 0s 26ms/step

	text	actual_labels	predicted_labels
0	rt	O	o
1	miriamsaying	O	o
2	may	O	o
3	mga	O	o
4	patama	O	o
5	talagang	O	o
6	di	O	o
7	naman	O	o
8	para	O	o
9	sayo	O	o
10	sadyang	O	o
11	assumera	O	o
12	ka	O	o
13	lang	O	o
14	kaya	O	o
15	inaangkin	O	o
16	mo	O	o

## BERT

```
In [65]: 1 # imports
          2 from transformers import TFBertForTokenClassification
          3 from transformers import BertTokenizer
```

```
In [66]: 1 # initializing the bert tokenizer (Sub word tokenizer)
          2 bert_tokenizer = BertTokenizer.from_pretrained('bert-base-uncased', do_low
```

```
In [67]: 1 # preprocessing data for bert
2
3 def preprocess_bert(text,labels):
4
5     text_list = []
6     label_list = []
7
8     for i in range(len(text)):
9         tokenized_words = []
10        tokenized_labels = []
11        words = text[i].split()
12        entities = labels[i].split()
13        sentence_len = len(words)
14        for j in range(sentence_len):
15            tokenized_word = bert_tokenizer.tokenize(words[j])
16            tokenized_words.extend(tokenized_word)
17            tokenized_label = [entities[j]]*len(tokenized_word)
18            tokenized_labels.extend(tokenized_label)
19            text_list.append(' '.join(tokenized_words))
20            label_list.append(' '.join(tokenized_labels))
21
22    return text_list,label_list
```

```
In [68]: 1 # applying preprocessing to the sentences
2 train_bert_tokenized_sentences, train_bert_tokenized_labels = preprocess_b
3 test_bert_tokenized_sentences, test_bert_tokenized_labels = preprocess_ber
```

```
In [69]: 1 # checking the sentences after tokenization
2 print(train_bert_tokenized_sentences[0])
3 print(train_bert_tokenized_labels[0])
```

new tee ##shi ##rts ordered with all new designs i hope people will like october 9th 2010 remington s annual october ##fest  
0 B-other I-other I-other

```
In [70]: 1 # adapt the vectorizer layer before modeling
2 sentence_indexer_bert = TextVectorization(max_tokens=NUM_WORDS, output_seq
3 sentence_indexer_bert.adapt(train_bert_tokenized_sentences)
4 train_bert_sentence_index = sentence_indexer_bert(train_bert_tokenized_se
5 test_bert_sentence_index = sentence_indexer_bert(test_bert_tokenized_sen
```

```
In [71]: 1 # initializing the token type ids for train and test sentences
2 train_token_type_ids = np.zeros(shape=(len(train_bert_tokenized_sentences),
3 test_token_type_ids = np.zeros(shape=(len(test_bert_tokenized_sentences),S
```

```
In [72]: 1 # initialiazing the attention masks for train and test sentences
2 train_attn_mask = np.zeros(shape=(len(train_bert_tokenized_sentences), SENTENCE_LENGTH))
3 test_attn_mask = np.zeros(shape=(len(test_bert_tokenized_sentences), SENTENCE_LENGTH))
4 for i in range(len(train_bert_tokenized_sentences)):
5     length = min(len(train_bert_tokenized_sentences[i].split()), SENTENCE_LENGTH)
6     train_attn_mask[i,:length] = 1
7 for i in range(len(test_bert_tokenized_sentences)):
8     length = min(len(test_bert_tokenized_sentences[i].split()), SENTENCE_LENGTH)
9     test_attn_mask[i,:length] = 1
```

```
In [73]: 1 # prepare the entity data
2
3 # tokenizing the entities
4 entity_indexer_bert = TextVectorization(max_tokens=NUM_ENTITIES, output_sequence_length=SENTENCE_LENGTH)
5 entity_indexer_bert.adapt(ytrain)
6 train_bert_label_index = entity_indexer_bert(train_bert_tokenized_labels)
7 test_bert_label_index = entity_indexer_bert(test_bert_tokenized_labels)
```

```
In [74]: 1 # imports for modeling
2 from keras.losses import SparseCategoricalCrossentropy
3 from keras.optimizers import Adam
```

In [87]:

```
1 # build bert model
2 encoder = TFBertForTokenClassification.from_pretrained('bert-base-uncased')
3 def build_bert_model():
4
5     # getting inputs to the model
6     input_sentence_ids = Input(shape=(SENTENCE_LENGTH,), dtype=tf.int32)
7     input_token_type_ids = Input(shape=(SENTENCE_LENGTH,), dtype=tf.int32)
8     input_attn_ids = Input(shape=(SENTENCE_LENGTH,), dtype=tf.int32)
9
10    # sending inputs through the bert model
11    embeddings = encoder(
12        input_ids = input_sentence_ids,
13        token_type_ids = input_token_type_ids,
14        attention_mask = input_attn_ids,
15    )[0]
16
17    # sending the context vectors through dense Layer with Linear activation
18    output_logits = Dense(NUM_ENTITIES, activation='linear')(embeddings)
19
20    # defining the model
21    model = Model(
22        inputs = [input_sentence_ids, input_token_type_ids, input_attn_ids],
23        outputs = [output_logits]
24    )
25
26    # compiling the model
27    model.compile(
28        loss = SparseCategoricalCrossentropy(from_logits=True),
29        optimizer = Adam(),
30        run_eagerly=True
31    )
32
33    # returning the built model
34    return model
```

All PyTorch model weights were used when initializing TFBertForTokenClassification.

Some weights or buffers of the TF 2.0 model TFBertForTokenClassification were not initialized from the PyTorch model and are newly initialized: ['classifier.weight', 'classifier.bias']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
In [88]: 1 # calling the build model function and checking it's structure
2 bert_model = build_bert_model()
3 bert_model.summary()
```

Model: "model\_3"

---

Layer (type) to	Output Shape	Param #	Connected to
<hr/>			
input_4 (InputLayer)	[(None, 45)]	0	[]
input_6 (InputLayer)	[(None, 45)]	0	[]
input_5 (InputLayer)	[(None, 45)]	0	[]
bert_layer (TFBertForToken Classification)	TFTokenClassifierOutput(lo ss=None, logits=(None, 45, 2), hidden_states=None, atten tions=None)	1088931	['input_4 [0][0]', 'input_6 [0][0]', 'input_5 [0][0]']
dense_3 (Dense)	(None, 45, 21)	63	['bert_la yer[0][0]']
<hr/>			
<hr/>			
Total params: 108893249 (415.39 MB)			
Trainable params: 108893249 (415.39 MB)			
Non-trainable params: 0 (0.00 Byte)			

---

```
In [89]: 1 # defining callbacks
2 mc = tf.keras.callbacks.ModelCheckpoint(
3     os.path.join(root_path, 'models', 'bert.ckpt'),
4     monitor="val_loss",
5     save_best_only=True,
6 )
7 es = tf.keras.callbacks.EarlyStopping(
8     monitor="val_loss",
9     min_delta=0,
10    patience=5,
11    start_from_epoch=5,
12 )
13 callbacks=[mc,es]
```

```
In [92]: 1 # fitting the bert model
2 history_bert = bert_model.fit(
3     [train_bert_sentence_index, train_token_type_ids, train_attn_mask], tr
4     validation_data = ([test_bert_sentence_index, test_token_type_ids, tes
5     epochs=3,
6     batch_size=16,
7     callbacks=callbacks
8 )
```

Epoch 1/3  
313/313 [=====] - ETA: 0s - loss: 0.9673INFO:tensorflow:Assets written to: c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\models\bert.ckpt\assets  
INFO:tensorflow:Assets written to: c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\models\bert.ckpt\assets  
313/313 [=====] - 1137s 4s/step - loss: 0.9673 - val\_loss: 0.9566  
Epoch 2/3  
313/313 [=====] - ETA: 0s - loss: 0.9585INFO:tensorflow:Assets written to: c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\models\bert.ckpt\assets  
INFO:tensorflow:Assets written to: c:\Users\psyki\Downloads\Learning\github-repos\twitter-ner-case-study\models\bert.ckpt\assets  
313/313 [=====] - 1240s 4s/step - loss: 0.9585 - val\_loss: 0.9535  
Epoch 3/3  
313/313 [=====] - 1303s 4s/step - loss: 0.9572 - val\_loss: 0.9539

```
In [93]: 1 # Loading the best model
2 bert_model.load_weights(os.path.join(root_path, 'models', 'bert.ckpt'))
```

Out[93]: <tensorflow.python.checkpoint.checkpoint.CheckpointLoadStatus at 0x26cbe72a650>

## BERT Inference

```
In [94]: 1 # getting the predictions
2 predictions = bert_model.predict([test_bert_sentence_index,test_token_type
```

40/40 [=====] - 56s 1s/step

In [95]:

```

1 # defining a function to join back the tokenized text
2 def get_joined_labels(tokenized_text,tokenized_labels):
3     joined_labels = []
4     begin = 0
5     end = 0
6
7     for i in range(len(tokenized_text.split())):
8
9         if i==len(tokenized_text.split())-1:
10             if begin==end:
11                 joined_labels.append(tokenized_labels.split()[i])
12             else:
13                 label_segment = tokenized_labels.split()[begin:end+1]
14                 added = None
15                 for j in label_segment:
16                     if j!='o':
17                         added = j
18                         break
19                 if added == None:
20                     added = 'o'
21                 joined_labels.append(added)
22
23             elif str(tokenized_text.split()[i+1]).startswith("##"):
24                 end = i+1
25
26             else:
27                 if begin==end:
28                     begin = i+1
29                     end = i+1
30                     joined_labels.append(tokenized_labels.split()[i])
31                 else:
32                     label_segment = tokenized_labels.split()[begin:end+1]
33                     added = None
34                     for j in label_segment:
35                         if j!='o':
36                             added = j
37                             break
38                     if added == None:
39                         added = 'o'
40                     joined_labels.append(added)
41                     begin = i+1
42                     end = i+1
43             return ''.join(joined_labels)

```

In [98]:

```

1 # joining back the tokenized text
2 pred_labels = np.argmax(predictions, axis=-1)
3 bert_joined_labels = []
4 for i in range(len(xtest)):
5     text = test_bert_tokenized_sentences[i]
6     lbls = ' '.join(np.asarray(entity_indexer_bert.get_vocabulary())[list(
7         bert_joined_labels.append(get_joined_labels(text,lbls)))

```

```
In [ ]: 1 # function to infer prediction from a text in xtest randomly
2 def infer_bert():
3     idx = np.random.choice(range(len(test_bert_sentence_index)), 1, replace=
4
5     text = xtest[idx]
6     labels_act = ytest[idx]
7     pred_labels = bert_joined_labels[idx]
8     result = pd.DataFrame(
9         {
10            'text': text.split(),
11            'actual_labels': labels_act.split(),
12            'predicted_labels': pred_labels.split()
13        }
14    )
15    display(result)
```

```
In [ ]: 1 # infer random text from the validation data
2 infer_bert()
```

1/1 [=====] - 0s 26ms/step

	text	actual_labels	predicted_labels
0	rt	O	o
1	miriamsaying	O	o
2	may	O	o
3	mga	O	o
4	patama	O	o
5	talagang	O	o
6	di	O	o
7	naman	O	o
8	para	O	o
9	sayo	O	o
10	sadyang	O	o
11	assumera	O	o
12	ka	O	o
13	lang	O	o
14	kaya	O	o
15	inaangkin	O	o
16	mo	O	o

## Conclusion

- The BI-LSTM+CRF model tends to overfit to the data
- The prediction from the CRF layer is not accurate

- Since BERT is larger model with model weights and more layers, it takes a lot of time to train
- The loss in BERT is larger compared to that in LSTM
- BERT tends to perform better than the LSRM model even though the loss is high
- Early stopping and Model checkpoint callbacks assist in training the model

## Questions

**Q.** Defining the problem statements and where can this and modifications of this be used?

**Ans.** Similar model can be used in Part-of-Speech tagging (POS tagging) or any other problem where we need a prediction at a token level. Meaning that the input length and the output length is the same.

**Q.** Explain the data format (conll bio format)

**Ans.** Conll format is a text storing format where each word in the text is separated by a line (\n), and annotation of that word is separated by tab (\t) and each text is separated by two lines (\n\n).

**Q.** What other ner data annotation formats are available and how are they different

**Ans.** Other formats can be BIO, IOB, JSON and XML

**Q.** Why do we need tokenization of the data in our case

**Ans.** We need tokenization in our case to divide the sentence into smaller constituents in order to capture the sequence, context and attention mechanism. If we do not tokenize our text, the input to the model will be one big string. Another problem will be to convert the string to vectors.

**Q.** What other models can you use for this task

**Ans.** We can replace LSTM with RNN or GRU and couple them with CRF layer. Or in case of BERT, we can use transformers (encoder-decoder architecture) or GPT architecture.

**Q.** Did early stopping have any effect on the training and results.

**Ans.** Yes. Since the model started to overfit as the epochs progresses, early stopping caused the training to stop since the performance was not improving.

**Q.** How does the BERT model expect a pair of sentences to be processed?

**Ans.** The BERT model expects a token\_type\_id tensor, which represents which sentence does each of the token belongs to. The pair of sentences are merged into a single vector and token\_type\_id is used to distinguish between them.

**Q.** Why choose Attention based models over Recurrent based ones?

**Ans.** Recurrent based models have problems capturing long term dependencies of the words. Also, each word can be connected to multiple other words which can be captured by the attention heads. Hence, attention based models are preferred if long term dependencies or connection of words is required.

**Q.** Differentiate BERT and simple transformers

**Ans.** BERT can be considered as the encoder only part of the transformers. In transformer architecture, first the input is passed through the encoder, then through the decoder and then the final prediction is made. But in BERT, the encoder only architecture is responsible for prediction.

In [ ]: 1

In [ ]: 1

In [ ]: 1