## **Text Generation Model using ML**

A Text Generation Model is a type of Natural Language Processing (NLP) model that automatically generates human-like text. It can produce coherent and contextually relevant text based on the input text. So, if you want to learn how to build a Text Generation Model, this article is for you. In this article, I'll take you through the task of building a Text Generation Model with Deep Learning using the Python programming language.

## **Text Generation Model:**

Text Generation Models have various applications, such as content creation, chatbots, automated story writing, and more. They often utilize advanced Machine Learning techniques, particularly Deep Learning models like Recurrent Neural Networks (RNNs), Long Short-Term Memory Networks (LSTMs), and Transformer models like GPT (Generative Pre-trained Transformer).

Below is the process we can follow for the task of building a Text Generation Model:

- Understand what you want to achieve with the text generation model (e.g., chatbot responses, creative writing, code generation).
- Consider the style, complexity, and length of the text to be generated.
- Collect a large dataset of text that's representative of the style and content you want to generate.
- Clean the text data (remove unwanted characters, correct spellings), and preprocess it (tokenization, lowercasing, removing stop words if necessary).
- Choose a deep neural network architecture to handle sequences for text generation.
- Frame the problem as a sequence modeling task where the model learns to predict the next words in a sequence.
- Use your text data to train the model.

For this task, we can use the Tiny Shakespeare dataset because of two reasons:

- It's available in the format of dialogues, so you will learn how to generate text in the form of dialogues.
- Usually, we need huge textual datasets for building text generation models. The Tiny Shakespeare dataset is already available in the tensorflow datasets, so we don't need to download any dataset externally.

So, let's understand how to build a Text Generation Model with Deep Learning using Python. I'll start this task by importing the necessary Python libraries and the dataset:

```
[1] 1 import tensorflow as tf
2 import tensorflow_datasets as tfds
3 import numpy as np
4
5 # load the Tiny Shakespeare dataset
6 dataset, info = tfds.load('tiny_shakespeare', with_info=True, as_supervised=False)

Downloading and preparing dataset Unknown size (download: Unknown size, generated: 1.06 MiB, total: 1.06 MiB) to /root/tensorflow_datasets/tiny_shakespeare/1.0.0...
DI Completed...: 100% 1/1 [00:00<00:00, 2.49 url/s]

DI Size...: 1/0 [00:00<00:00, 2.77 MiB/s]</pre>
```

Our dataset contains data in a textual format. Language models need numerical data, so we'll convert the text to sequences of integers. We'll also create sequences for training:

Dataset tiny\_shakespeare downloaded and prepared to /root/tensorflow\_datasets/tiny\_shakespeare/1.0.0. Subsequent calls will reuse this data.

```
1 # get the text from the dataset
2 text = next(iter(dataset['train']))['text'].numpy().decode('utf-8')
3
4 # create a mapping from unique characters to indices
5 vocab = sorted(set(text))
6 char2idx = {char: idx for idx, char in enumerate(vocab)}
7 idx2char = np.array(vocab)
8
9 # numerically represent the characters
10 text_as_int = np.array([char2idx[c] for c in text])
11
12 # create training examples and targets
13 seq_length = 100
14 examples_per_epoch = len(text) // (seq_length + 1)
15
16 # create training sequences
17 char_dataset = tf.data.Dataset.from_tensor_slices(text_as_int)
18
19 sequences = char_dataset.batch(seq_length + 1, drop_remainder=True)
```

For each sequence, we will now duplicate and shift it to form the input and target text by using the map method to apply a simple function to each batch:

```
1 def split_input_target(chunk):
2    input_text = chunk[:-1]
3    target_text = chunk[1:]
4    return input_text, target_text
5
6 dataset = sequences.map(split_input_target)
```

Now, we'll shuffle the dataset and pack it into training batches:

```
1 # batch size and buffer size
2 BATCH_SIZE = 64
3 BUFFER_SIZE = 10000
4
5 dataset = (
6    dataset
7    .shuffle(BUFFER_SIZE)
8    .batch(BATCH_SIZE, drop_remainder=True)
9    .prefetch(tf.data.experimental.AUTOTUNE)
10 )
```

Now, we'll use a simple Recurrent Neural Network model with a few layers to build the model:

```
Suggested code may be subject to a license | stackoverflow.com/questions/64844743/what-does-it-mean-when-the-loss-starts-going-up-again
 1 # length of the vocabulary
 2 vocab_size = len(vocab)
4 # the embedding dimension
 5 embedding_dim = 256
7 # number of RNN units
 8 rnn_units = 1024
9
10 def build_model(vocab_size, embedding_dim, rnn_units, batch_size):
11 model = tf.keras.Sequential([
         tf.keras.layers.Embedding(vocab_size, embedding_dim, batch_input_shape=[batch_size, None]),
       tf.keras.layers.LSTM(rnn_units, return_sequences=True, stateful=True, recurrent_initializer='glorot_uniform'), tf.keras.layers.Dense(vocab_size)
13
14
     ])
      return model
16
17 model = build_model(vocab_size,embedding_dim,rnn_units, BATCH_SIZE)
```

We'll now choose an optimizer and a loss function to compile the model:

```
1 def loss(labels, logits):
2    return tf.keras.losses.sparse_categorical_crossentropy(labels, logits, from_logits=True)
3
4 model.compile(optimizer='adam', loss=loss)
```

## We'll now train the model:

```
1 import os
2
3 # directory where the checkpoints will be saved
4 checkpoint_dir = './training_checkpoints'
5
6 # name of the checkpoint files
7 checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt_{epoch}")
8
9 checkpoint_callback = tf.keras.callbacks.ModelCheckpoint(
10 filepath=checkpoint_prefix,
11 save_weights_only=True
12 )
13
14 # train the model
15 EPOCHS = 10
16 history = model.fit(dataset, epochs=EPOCHS, callbacks=[checkpoint_callback])
```

```
→ Epoch 1/10

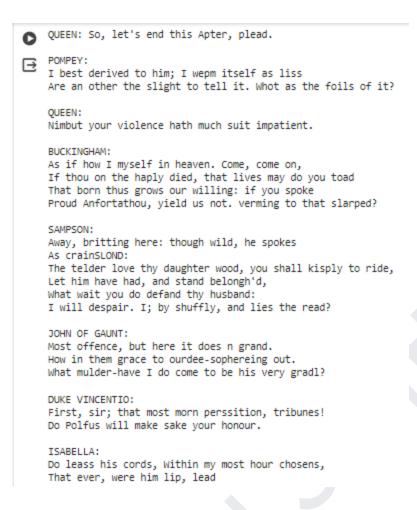
   155/155 [=============== ] - 1225s 8s/step - loss: 2.6506
   Epoch 2/10
   155/155 [============] - 1184s 8s/step - loss: 1.9349
   Epoch 3/10
   155/155 [============== ] - 1181s 8s/step - loss: 1.6805
   Epoch 4/10
   155/155 [============ ] - 1169s 8s/step - loss: 1.5409
   Epoch 5/10
   Epoch 6/10
   155/155 [============= ] - 1179s 8s/step - loss: 1.3941
   Epoch 7/10
   155/155 [============ ] - 1176s 8s/step - loss: 1.3473
   Epoch 8/10
   155/155 [============ ] - 1172s 8s/step - loss: 1.3089
   Epoch 9/10
   Epoch 10/10
   155/155 [=============== ] - 1194s 8s/step - loss: 1.2394
```

After training, we can now use the model to generate text. First, we will restore the latest checkpoint and rebuild the model with a batch size of 1:

```
[ ] 1 model = build_model(vocab_size, embedding_dim, rnn_units, batch_size=1)
2 model.load_weights(tf.train.latest_checkpoint(checkpoint_dir))
3 model.build(tf.TensorShape([1, None]))
```

Now, to generate text, we'll input a seed string, predict the next character, and then add it back to the input, continuing this process to generate longer text:

```
1 def generate_text(model, start_string):
      num_generate = 1000
3
4
      input_eval = [char2idx[s] for s in start_string]
5
      input eval = tf.expand dims(input eval, 0)
6
7
      text_generated = []
8
9
      model.reset_states()
10
      for i in range(num generate):
11
          predictions = model(input eval)
          predictions = tf.squeeze(predictions, 0)
12
13
          predicted_id = tf.random.categorical(predictions, num_samples=1)[-1, 0].numpy()
14
          input_eval = tf.expand_dims([predicted_id], 0)
15
16
17
          text generated.append(idx2char[predicted id])
18
      return (start_string + ''.join(text_generated))
19
20 print(generate_text(model, start_string = u"QUEEN: So, let's end this"))
```



The generate\_text function in the above code uses a trained Recurrent Neural Network model to generate a sequence of text, starting with a given seed phrase (start\_string). It converts the seed phrase into a sequence of numeric indices, feeds these indices into the model, and then iteratively generates new characters, each time using the model's most recent output as the input for the next step. This process continues for a specified number of iterations (num\_generate), resulting in a stream of text that extends from the initial seed.

The function employs randomness in character selection to ensure variability in the generated text, and the final output is a concatenation of the seed phrase with the newly generated characters, typically reflecting the style and content of the training data used for the model.

## **Summary**

So, this is how we build a Text Generation Model with Deep Learning using Python. Text Generation Models have various applications, such as content creation, chatbots, automated story writing, and more. They often utilize advanced Machine Learning techniques, particularly Deep Learning models like Recurrent Neural Networks (RNNs), Long Short-Term Memory Networks (LSTMs), and Transformer models like GPT (Generative Pre-trained Transformer).

Anmol Sonthalia GitHub - anmol957