

Final_Chatbot_Model2_Kaustubh_Prakash_Deolase.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text | Run all

RAM Disk

Final Project Chatbot Model by Kaustubh Prakash Deolase

ChatBot Model Seq2Seq Chatbot (LSTM) by Kaustubh Prakash Deolase

```
[16] ✓ 0s
1 import os
2 import re
3 import numpy as np
4 import pandas as pd
5 import tensorflow as tf
6 from tensorflow.keras.layers import (Embedding, LSTM, Dense, LayerNormalization, Dropout, TextVectorization)
7 import matplotlib.pyplot as plt
```

TEXT CLEANING FUNCTION

```
[17] ✓ 0s
1 def clean_text(text: str) -> str:
2     """Basic text cleaning."""
3     if not isinstance(text, str):
4         text = str(text)
5     text = text.lower()
6     text = re.sub(r"-", " ", text)
7     text = re.sub(r"\([.,!]\)", r"\1 ", text)
8     text = re.sub(r"\([0-9]\)", r"\1 ", text)
9     text = re.sub(r"\s+", " ", text).strip()
10    return text
```

LOAD AND PREPROCESS DIALOGS

```
[18] ✓ 0s
1 def load_dialogs(path="dialogs.txt"):
2     if not os.path.exists(path):
3         raise FileNotFoundError(f"dialogs file not found at: {path}")
4     df = pd.read_csv(path, sep="\t", names=["question", "answer"], encoding="utf-8")
5     df.dropna(inplace=True)
6     df["question"] = df["question"].map(clean_text)
7     df["answer"] = df["answer"].map(clean_text)
8     return df
```

TEXT VECTORIZATION

```
[19] ✓ 0s
1 def build_vectorizer(texts, max_tokens=5000, seq_len=30):
2     vectorizer = TextVectorization(
3         max_tokens=max_tokens,
4         output_mode="int",
5         output_sequence_length=seq_len,
6         standardize=None # We already cleaned text
7     )
8     vectorizer.adapt(texts)
9     return vectorizer
```

ENCODER MODEL

```
[20] ✓ 0s
1 class Encoder(tf.keras.Model):
2     def __init__(self, units, embedding_dim, vocab_size):
3         super().__init__()
4         self.embedding = Embedding(vocab_size, embedding_dim, mask_zero=True)
5         self.norm = LayerNormalization()
6         self.dropout = Dropout(0.4)
7         self.lstm = LSTM(units, return_sequences=True, return_state=True, dropout=0.4)
8
9     def call(self, x, training=False):
10        x = self.embedding(x)
11        x = self.norm(x, training=training)
12        x = self.dropout(x, training=training)
13        outputs, h, c = self.lstm(x, training=training)
14        return outputs, h, c
```

DECODER MODEL

```
[21] ✓ 0s
1 class Decoder(tf.keras.Model):
2     def init (self, units, embedding_dim, vocab_size):
```

```

3     super().__init__(self)
4     self.embedding = Embedding(vocab_size, embedding_dim, mask_zero=True)
5     self.norm1 = LayerNormalization()
6     self.lstm = LSTM(units, return_sequences=True, return_state=True, dropout=0.4)
7     self.norm2 = LayerNormalization()
8     self.dropout = Dropout(0.4)
9     self.fc = Dense(vocab_size, activation="softmax")
10
11    def call(self, x, initial_state=None, training=False):
12        x = self.embedding(x)
13        x = self.norm1(x, training=training)
14        x = self.dropout(x, training=training)
15        outputs, h, c = self.lstm(x, initial_state=initial_state, training=training)
16        outputs = self.norm2(outputs, training=training)
17        outputs = self.dropout(outputs, training=training)
18        logits = self.fc(outputs)
19        return logits, h, c

```

BUILD TRAINING MODEL

```

[22]
1 def build_training_model(encoder, decoder, seq_len):
2     enc_inp = tf.keras.Input(shape=(seq_len,), name="encoder_input")
3     dec_inp = tf.keras.Input(shape=(seq_len,), name="decoder_input")
4     enc_out, h, c = encoder(enc_inp)
5     dec_out, _, _ = decoder(dec_inp, initial_state=[h, c])
6     model = tf.keras.Model([enc_inp, dec_inp], dec_out)
7     return model

```

BUILD INFERENCE MODELS (for chatting)

```

[23]
1 def build_inference_models(encoder, decoder, seq_len, units):
2     enc_input = tf.keras.Input(shape=(seq_len,), name="enc_input")
3     _, h, c = encoder(enc_input)
4     enc_model = tf.keras.Model(enc_input, [h, c])
5
6     dec_input = tf.keras.Input(shape=(1,), name="dec_input")
7     h_in = tf.keras.Input(shape=(units,), name="h_in")
8     c_in = tf.keras.Input(shape=(units,), name="c_in")
9     dec_out, h_out, c_out = decoder(dec_input, initial_state=[h_in, c_in])
10    dec_model = tf.keras.Model([dec_input, h_in, c_in], [dec_out, h_out, c_out])
11    return enc_model, dec_model

```

GREEDY DECODING

```

[24]
1 def greedy_decode(input_text, vectorizer, enc_model, dec_model, id2word, word2id, max_len=20):
2     # Convert text to sequence using TextVectorization
3     seq = vectorizer([input_text]) # returns tensor (1, seq_len)
4
5     # Encode input sequence to get states
6     state_h, state_c = enc_model(seq)
7
8     # Start token id
9     start_id = word2id.get('<start>', 1)
10    end_id = word2id.get('<end>', 2)
11
12    target_seq = tf.expand_dims([start_id], 0) # shape (1,1)
13    decoded_sentence = []
14
15    for _ in range(max_len):
16        output_tokens, h, c = dec_model([target_seq, state_h, state_c])
17        sampled_token_index = tf.argmax(output_tokens[0, -1, :]).numpy()
18        sampled_word = id2word.get(sampled_token_index, '')
19
20        if sampled_word == '<end>':
21            break
22
23        decoded_sentence.append(sampled_word)
24        target_seq = tf.expand_dims([sampled_token_index], 0)
25        state_h, state_c = h, c
26
27    return ' '.join(decoded_sentence)
28
29

```

MAIN TRAINING PIPELINE

```

[25]
1 def train_chatbot(data_path="dialogs.txt"):
2     # Hyperparameters
3     VOCAB_SIZE = 5000
4     SEQ_LEN = 30
5     EMBED_DIM = 128
6     LSTM_UNITS = 256
7     EPOCHS = 30
8     BATCH_SIZE = 128
9     LR = 1e-3
10

```

```

11 # Load data
12 df = load_dialogs(data_path)
13 df["decoder_input"] = "<start> " + df["answer"]
14 df["decoder_target"] = df["answer"] + " <end>"
15
16 # Build vocabulary
17 combined_text = pd.concat([
18     df['question'], df['decoder_input'], df['decoder_target']
19 ])
20 vectorizer = build_vectorizer(combined_text, max_tokens=VOCAB_SIZE, seq_len=SEQ_LEN)
21
22 vocab = vectorizer.get_vocabulary()
23 word2id = {t: i for i, t in enumerate(vocab)}
24 id2word = {i: t for i, t in enumerate(vocab)}
25
26 # Vectorize sequences
27 enc_inp = vectorizer(df["question"])
28 dec_inp = vectorizer(df["decoder_input"])
29 dec_tar = vectorizer(df["decoder_target"])
30
31 # Build model
32 encoder = Encoder(LSTM_UNITS, EMBED_DIM, VOCAB_SIZE)
33 decoder = Decoder(LSTM_UNITS, EMBED_DIM, VOCAB_SIZE)
34 model = build_training_model(encoder, decoder, SEQ_LEN)
35 model.compile(
36     optimizer=tf.keras.optimizers.Adam(LR),
37     loss=tf.keras.losses.SparseCategoricalCrossentropy(),
38     metrics=["accuracy"]
39 )
40
41 # Train
42 ckpt = tf.keras.callbacks.ModelCheckpoint(
43     "models/best_seq2seq.keras", monitor="val_loss", save_best_only=True
44 )
45
46 history = model.fit(
47     [enc_inp, dec_inp], dec_tar,
48     batch_size=BATCH_SIZE,
49     epochs=EPOCHS,
50     validation_split=0.1,
51     callbacks=[ckpt]
52 )
53
54 # Save models
55 enc_model, dec_model = build_inference_models(encoder, decoder, SEQ_LEN, LSTM_UNITS)
56 os.makedirs("models", exist_ok=True)
57 model.save("models/seq2seq_model.keras")
58 enc_model.save("models/encoder_model.keras")
59 dec_model.save("models/decoder_model.keras")
60
61 # Plot metrics
62 fig, ax = plt.subplots(1, 2, figsize=(15, 5))
63 ax[0].plot(history.history["loss"], label="train_loss", color="red")
64 ax[0].plot(history.history["val_loss"], label="val_loss", color="blue")
65 ax[0].set_title("Loss")
66 ax[0].legend()
67 ax[1].plot(history.history["accuracy"], label="train_acc", color="green")
68 ax[1].plot(history.history["val_accuracy"], label="val_acc", color="orange")
69 ax[1].set_title("Accuracy")
70 ax[1].legend()
71 plt.show()
72
73 print("✅ Training complete and models saved in 'models/' folder.")
74 return vectorizer, enc_model, dec_model, id2word, word2id

```

CHAT LOOP (AFTER TRAINING)

```

[26] 1 def chat_loop(vectorizer, enc_model, dec_model, id2word, word2id):
2     print("\n🤖 Chatbot is ready! Type 'exit' to quit.\n")
3     while True:
4         msg = input("You: ").strip()
5         if msg.lower() in ["exit", "quit", "bye"]:
6             print("Chatbot: Goodbye!")
7             break
8         reply = greedy_decode(clean_text(msg), vectorizer, enc_model, dec_model, id2word, word2id)
9         print("Chatbot:", reply)

```

ENTRY POINT

```

[27] 1 if __name__ == "__main__":
2     vectorizer, enc_model, dec_model, id2word, word2id = train_chatbot("dialogs.txt")
3     #chat_loop(vectorizer, enc_model, dec_model, id2word, word2id)
4
5 Epoch 0/30
6 2s 57ms/step - accuracy: 0.1151 - loss: 2.9962 - val_accuracy: 0.0897 - val_loss: 4.8993
7 Epoch 1/30
8 2s 57ms/step - accuracy: 0.1187 - loss: 2.9018 - val_accuracy: 0.0887 - val_loss: 4.9450
9 Epoch 2/30
10 2s 57ms/step - accuracy: 0.1195 - loss: 2.8625 - val_accuracy: 0.0896 - val_loss: 4.9812
11 Epoch 3/30
12 2s 57ms/step - accuracy: 0.1224 - loss: 2.7705 - val_accuracy: 0.0871 - val_loss: 5.0308
13 Epoch 4/30
14 2s 59ms/step - accuracy: 0.1249 - loss: 2.7118 - val_accuracy: 0.0874 - val_loss: 5.0600
15 Epoch 5/30
16 2s 62ms/step - accuracy: 0.1293 - loss: 2.6244 - val_accuracy: 0.0885 - val_loss: 5.1169
17 Epoch 6/30
18 2s 61ms/step - accuracy: 0.1335 - loss: 2.5278 - val_accuracy: 0.0859 - val_loss: 5.1414
19 Epoch 7/30
20 2s 56ms/step - accuracy: 0.1336 - loss: 2.4935 - val_accuracy: 0.0869 - val_loss: 5.1859
21 Epoch 8/30
22 2s 56ms/step - accuracy: 0.1378 - loss: 2.4160 - val_accuracy: 0.0874 - val_loss: 5.2255
23 Epoch 9/30
24 2s 56ms/step - accuracy: 0.1406 - loss: 2.3588 - val_accuracy: 0.0861 - val_loss: 5.2673
25 Epoch 10/30
26 2s 57ms/step - accuracy: 0.1425 - loss: 2.2994 - val_accuracy: 0.0850 - val_loss: 5.3048

```

27/27 2s 57ms/step - accuracy: 0.1428 - loss: 2.2534 - val_accuracy: 0.0842 - val_loss: 5.3378
 Epoch 27/30
 27/27 2s 57ms/step - accuracy: 0.1462 - loss: 2.1978 - val_accuracy: 0.0849 - val_loss: 5.3791
 Epoch 28/30
 27/27 2s 65ms/step - accuracy: 0.1474 - loss: 2.1679 - val_accuracy: 0.0858 - val_loss: 5.4156
 Epoch 29/30
 27/27 2s 59ms/step - accuracy: 0.1528 - loss: 2.0927 - val_accuracy: 0.0865 - val_loss: 5.4488
 Epoch 30/30
 27/27 2s 56ms/step - accuracy: 0.1556 - loss: 2.0482 - val_accuracy: 0.0855 - val_loss: 5.4963

The figure consists of two line plots side-by-side. The left plot, titled 'Loss', shows training loss (red line) decreasing from approximately 6.2 at epoch 0 to about 2.2 at epoch 30, while validation loss (blue line) remains relatively flat, starting around 5.4 and ending near 5.5. The right plot, titled 'Accuracy', shows training accuracy (green line) increasing from 0.05 at epoch 0 to 0.13 at epoch 30, and validation accuracy (orange line) increasing from 0.07 at epoch 0 to 0.13 at epoch 30.

ChatBot Conversation:

```
[28] ✓ 0s
1 # Function to reply to a single user input
2 def chatbot_reply(user_input, vectorizer, enc_model, dec_model, id2word, word2id):
3     user_input = clean_text(user_input)
4     return greedy_decode(user_input, vectorizer, enc_model, dec_model, id2word, word2id)
5
6 # Function to simulate a conversation from a list of messages
7 def print_conversation(messages, vectorizer, enc_model, dec_model, id2word, word2id):
8     for text in messages:
9         print(f'You: {text}')
10        reply = chatbot_reply(text, vectorizer, enc_model, dec_model, id2word, word2id)
11        print(f'Bot: {reply}')
12        print('=====')
```

You: hi
Bot: i'm not sleeping .
=====

You: how are you?
Bot: i don't know .
=====

You: what is your name?
Bot: i don't know .
=====

You: tell me a joke
Bot: i hope it doesn't rain
=====

You: bye
Bot: not yet .
=====