

WORKING CHATBOT IN PYTHON

PROJECT REPORT

INTRODUCTION TO MACHINE LEARNING

By

PARAS WAHI E18ECE026

Submitted to

Dr. Manoj Sharma



INDEX

1. INTRODUCTION TO CHATBOT
2. SETTING RULES FOR PRE DEFINED TRANSACTIONAL RESPONSE
3. INTRODUCTION OF WORD EMBEDDINGS AND TF-IDF VECTORIZATION
4. MAJOR LANGUAGE MODELS AND THEIR USES
5. SENTIMENT ANALYSIS
6. SUMMARY

INTRODUCTION TO CHATBOT AND THE TECHNOLOGY BEHIND IT

Chatbots utilize the technology colloquially known as natural language processing, or NLP. NLP can be classified into two basic approaches. These approaches are rule based NLP and Statistical NLP, or Deep Learning based NLP.

Rule based NLP involves manually engineering a system of grammar and syntax rules. The rule based methods are becoming less popular today.

There are certain advantages. For instance, sentence level tasks, especially in the limited but well defined domain, can be excellent tools for parsing and extraction. This makes this approach ideal for query analysis or simple specific casts.

The reason rule based techniques will continue to be used in combination with other techniques is because of its ability to capture specific language, phenomenon and code. Those rules directly into the syntax during parsing rule based NLP are precise.

One historical example of rule based NLP is the old style text adventure games from the 19 eighties, such as Zork. These primitive systems inferred intention based on type descriptions of actions that the player wanted the character to take. Some disadvantages to rule based NLP include the requirement for manual rule development, slow person speed due to branching code and the inability to cover the full scope of human language when it comes to chatting back. Specifically, they're often referred to as transactional chatbots in general, and generally I represent the list of options to the customer or simply detect predetermined intentions within well defined limitations. The other method involves using machine learning techniques to train a machine learning models, such as a neural network, on a large data set of relevant text data. There are many advantages to this approach, and those advantages become more obvious with larger data sets and models with larger feature sets.

This becomes one of the main drawbacks of the statistical approach.

If you want it done right, you need to trade a very large model on a very large amount of data.

However, when done right, the results could be amazing. The famous GPT 3 natural language model from Open A I, for instance, is based on machine learning approach, which is utilizing a type of model called the Transformer a couple disadvantages of this method, or that it can be difficult to debug since machine learning models are black box entities, meaning they learn on their own from the training data and interpreting how they learned is not possible without extensive probing and testing through external inputs and outputs.

Another issue is that the model has no built in understanding of the language. It is trained on meaning that, especially with smaller data sets specific language phenomena might not be picked up by the training and lead The might lead to misinterpretation and syntactical error. These downsides come partly or mostly be mitigated first by increasing the size of the data set and tuning the model parameters.

And second, by combining methods from the rule based approach to the statistical approach, There are a number of third party libraries that are of enormous value to the process of building a chat bot.

SETTING RULES FOR PRE-DEFINED TRANSACTIONAL RESPONSES

```
1 import nltk
2 from nltk.chat.util import Chat, reflections
3
4 from io import StringIO
5 import sys
6
7 set_pairs = [
8     [
9         r"my name is (.*)",
10        ["Hello %1, How are you doing today ?"],
11    ],
12    [
13        r"hi|hey|hello",
14        ["Hello", "Hey there"],
15    ],
16    [
17        r"what is your name?",
18        ["You can call me a chatbot ?"],
19    ],
20    [
21        r"how are you ?",
22        ["I am fine, thank you! How can i help you?"],
23    ],
24    [
25        r"I am fine, thank you",
26        ["great to hear that, how can i help you?"],
27    ],
28
37    r"(.*) thank you so much, that was helpful",
38    ["I am happy to help", "No problem, you're welcome"],
39 ],
40 [
41     r"quit",
42     ["Bye, take care. See you soon :) ", "It was nice talking to you. See you soon :)"]
43 ],
44 ]
45
46
47
48
49 def chatbot():
50     print("Hi, I'm your automated assistant")
51     chat = Chat(set_pairs, reflections)
52     chat.converse():
53     old_stdout = sys.stdout
54     sys.stdout = mystdout = StringIO()
55     sys.stdout = old_stdout
56     print('farts' mystdout.getvalue())
57
58 def converse():
59     if
60
61 if __name__ == "__main__":
62     chatbot()
63
```

INTRODUCTION OF WORD EMBEDDINGS AND TF-IDF VECTORIZATION AND COSINE SIMILARITY

Importing necessary libraries

```
1 import nltk
2 import numpy as np
3 import random
4 import string # to process standard python strings
5 from nltk.chat.util import Chat, reflections
6 from sklearn.feature_extraction.text import TfidfVectorizer
7 from sklearn.metrics.pairwise import cosine_similarity
8
9 from io import StringIO
10 import sys
11 import json
12
13
14 import os
15
16 import requests
17 from requests import get
18 from lxml import html
19 from bs4 import BeautifulSoup
20
21 import web_query as wq
22 import combine_pdf
```

Pre-defined set pairs

```

25 flag = True
26
27 set_pairs = [
28     [
29         r"my name is (.*)",
30         ["Hello %1, How are you doing today ?",]
31     ],
32     [
33         r"hi|hey|hello",
34         ["Hello", "Hey there",]
35     ],
36     [
37         r"what is your name?",
38         ["You can call me a chatbot ?",]
39     ],
40     [
41         r"how are you ?",
42         ["I am fine, thank you! How can i help you?",]
43     ],
44     [
45         r"I am fine, thank you",
46         ["great to hear that, how can i help you?",]
47     ],
48     [
49         r"how can i help you? ",
50         ["It is my function to help you.", "I just want to fulfill my purpose",]
51     ],

```

Reading file and downloading databases

```

67 f=open(r'C:\Users\darf3\Documents\FLG Work\Rhyne Python Chore Bot\corpus.txt','r',errors = 'ignore')
68 ##raw=json.load(f)
69 ##raw= json.dumps(raw)
70 ##seperator = ' '
71 ##final_str = seperator.join(raw)
72
73 raw=f.read()
74
75 raw=raw.lower()# converts to lowercase
76 nltk.download('punkt') # first-time use only
77 nltk.download('wordnet') # first-time use only
78 sent_tokens = nltk.sent_tokenize(raw)# converts to list of sentences
79 word_tokens = nltk.word_tokenize(raw)# converts to list of words
80
81
82 lemmer = nltk.stem.WordNetLemmatizer()
83 #WordNet is a semantically-oriented dictionary of English included in NLTK.
84 def LemTokens(tokens):
85     return [lemmer.lemmatize(token) for token in tokens]
86 remove_punct_dict = dict((ord(punct), None) for punct in string.punctuation)
87 def LemNormalize(text):
88     return LemTokens(nltk.word_tokenize(text.lower().translate(remove_punct_dict)))
89
90
91 GREETING_INPUTS = ("hello", "hi", "greetings", "sup", "what's up","hey",)
92 GREETING_RESPONSES = ["hi", "hey", "*nods*", "hi there", "hello", "I am glad! You are talking to me"]
93

```

Defining the functions

```

100 def greeting(sentence):
101
102     for word in sentence.split():
103         if word.lower() in GREETING_INPUTS:
104             return random.choice(GREETING_RESPONSES)
105
106
107 def chatbot():
108
109     chat = Chat(set_pairs, reflections)
110
111     user_input = quit
112     try:
113         user_input = input(">")
114     except EOFError:
115         print(user_input)
116     if user_input:
117
118         user_input = user_input[:-1]
119         if chat.respond(user_input) != None:
120             print(chat.respond(user_input))
121         else:
122             user_response = user_input
123             user_response=user_response.lower()
124             if(user_response!='bye'):
125                 if(user_response=='thanks' or user_response=='thank you'):
126                     flag=False
127
128             print("Bot: You are welcome..")
129         else:
130             if(greeting(user_response)!= None):
131                 print("Bot: "+greeting(user_response))
132             else:
133
134                 if("python" in user_response):
135                     print("Bot: ",end="")
136                     print(response(user_response))
137                     sent_tokens.remove(user_response)
138                 elif("combine" and "file" in user_response):
139                     for i in range(100):
140                         print ("\n")
141                         print("Entering file combining mode. Please enter the exact directory we will be
142 combining.")
143
144                     directory = input(">")
145                     os.chdir(directory)
146                     print("Thank you. Now please enter keyword to identify target files by title:")
147                     keyword = input(">")
148                     reg_pattern = r'(*'+keyword+'.pdf)|(^'+keyword+'_.*\pdf)|(. *'+keyword+'\.pdf)|
149 (*'+keyword+'.pdf)'
150
151                     # instantiate a bot object, I call it bot1 here
152                     bot1 = combine_pdf.FileBots(directory, reg_pattern)
153                     # use the .locate method to find files of interest
154                     bot1.locate(show=True)
155                     bot1.pdf_merge_tree()
156                 else:

```

```

155                                     flag=False
156                                     print("Bot: Bye! take care..")
157     #reply = chat.respond()
158     ##         old_stdout = sys.stdout
159     ##         sys.stdout = mystdout = StringIO()
160     ##         sys.stdout = old_stdout
161     ##         print('farts',mystdout.getvalue())
162
163
164
165
166
167
168 def response(user_response):
169     robo_response=''
170     sent_tokens.append(user_response)
171     TfIdfVec = TfIdfVectorizer(tokenizer=LemNormalize)
172     tfidf = TfIdfVec.fit_transform(sent_tokens)
173     vals = cosine_similarity(tfidf[-1], tfidf)
174     idx=vals.argsort()[0][-2]
175     flat = vals.flatten()
176     flat.sort()
177     req_tfidf = flat[-2]
178     if(req_tfidf==0):
179
180         return None
181
182     else:
183         robo_response = robo_response+sent_tokens[idx]
184         return robo_response
185
186
187 if __name__ == "__main__":
188     start=True
189
190     while flag == True:
191         if start==True:
192             print("Hi, I'm your automated assistant")
193             start=False
194             chatbot()
195
196
197
198

```


MAJOR LANGUAGE MODELS AND THEIR USES

Generative pre-trained Transformer 3 or G. P. T. 3 is an auto progressive language model that uses deep learning to produce human-like text. It is the third generation language prediction model in the GPT series, created by Open A.I, a San Francisco based artificial intelligence research laboratory.

31 open A.I researchers and engineers presented the original May 28th 2020 paper introducing GPT 3.

The team increased the capacity of GPT 3 by over two orders of magnitude from that of its predecessor, GPT 2 making GPT 3, the largest non sparse

language model to date. GPT 3 is higher number of parameters granted a higher level accuracy relative to previous versions with smaller

capacity. GPT 3s capacity is 10 times larger than that of Microsoft's Turing NLG.

In his July 29th 2020 review in The New York Times, Farhad

Manjoo said that GPT 3, which can generate computer code and poetry as well as prose, is not just amazing spooky and humbling, but also more than a little terrifying.

Australian philosopher David Chalmers described GPT 3 as one of the most interesting and important AI systems ever produced, a review on Wired said that GPT 3 was provoking chills across Silicon Valley.

The National Law Review said that GPT 3 is an impressive step in the larger process, with Open A I and others finding useful applications for all this power while continuing

to work toward a more general intelligence, the G P T 3 model is huge, too large for the average person to train at home Open A. I does have a pay scale for access

to the GPT 3 API.

SENTIMENT ANALYSIS AND INTENT CLASSIFICATION

```
1 import os
2 import model, sample, encoder
3
4 import json
5 import os
6 import numpy as np
7 import tensorflow as tf
8 import gpt_2_simple as gpt2
9
10 model_name = "345M" # The GPT-2 model we're using
11
12 gpt2.download_gpt2(model_name=model_name) # Download the model
13
14
15 sess = gpt2.start_tf_sess()
16 gpt2.finetune(sess,
17 file_name,
18 model_name=model_name,
19 steps=400) # steps is max number of training steps
20
21 def interact_model(
22     model_name,
23     seed,
24     nsamples,
25     batch_size,
26     length,
27     temperature,
28
29     top_k,
30     models_dir
31 ):
32     models_dir = os.path.expanduser(os.path.expandvars(models_dir))
33     if batch_size is None:
34         batch_size = 1
35     assert nsamples % batch_size == 0
36
37     enc = encoder.get_encoder(model_name, models_dir)
38     hparams = model.default_hparams()
39     with open(os.path.join(models_dir, model_name, 'hparams.json')) as f:
40         hparams.override_from_dict(json.load(f))
41
42     if length is None:
43         length = hparams.n_ctx // 2
44     elif length > hparams.n_ctx:
45         raise ValueError("Can't get samples longer than window size: %s" % hparams.n_ctx)
46
47     with tf.Session(graph=tf.Graph()) as sess:
48         context = tf.placeholder(tf.int32, [batch_size, None])
49         np.random.seed(seed)
50         tf.set_random_seed(seed)
51         output = sample.sample_sequence(
52             hparams=hparams, length=length,
53             context=context,
54             batch_size=batch_size,
55             temperature=temperature, top_k=top_k
```

```

54         temperature=temperature, top_k=top_k
55     )
56
57     saver = tf.train.Saver()
58     ckpt = tf.train.latest_checkpoint(os.path.join(models_dir, model_name))
59     saver.restore(sess, ckpt)
60
61     while True:
62         raw_text = input("Model prompt >>> ")
63         while not raw_text:
64             print('Prompt should not be empty!')
65             raw_text = input("Model prompt >>> ")
66         context_tokens = enc.encode(raw_text)
67         generated = 0
68         for _ in range(nsamples // batch_size):
69             out = sess.run(output, feed_dict={
70                 context: [context_tokens for _ in range(batch_size)]
71             })[:, len(context_tokens):]
72             for i in range(batch_size):
73                 generated += 1
74                 text = enc.decode(out[i])
75                 print("=" * 40 + " SAMPLE " + str(generated) + " " + "=" * 40)
76                 print(text)
77         print("=" * 80)
78

```

SUMMARY

We can combine the above methods to create an adaptable and task oriented chatbot. It would be able to answer complex questions and also after training on a larger dataset be able to answer a variety of questions and also give out human-like responses.