# Machine Learning Assignment

## Machine Learning for Cyber Security

BSc (Hons.) in Information Technology Specializing in Cyber Security

Sri Lanka Institute of Information Technology

Sri Lanka

May 2020

IT17354516

B I Sariffodeen

# Contents

# 1. Introduction

## 1.1. Objective

Construct a simple, functional chatbot using a Machine Learning module.

## 1.2. Prerequisites

- Anaconda - Free and Open Source Distribution of Python Language to simplify Package Management and distribution
- Jupyter Notebook on conda - Open Source Web Application to develop live code
- TensorFlow - Open Source software library for data flow and differentiable programming
- Keras - Open Source Neural Network library in python
- NLTK - Symbolic and Statistical Natural Language Processing
- nvidia CUDA - Parallel Computing Platform and API model
- Pickle - Pickling converts objects hierarchy to byte stream and unpickling does the reverse
- Azure and Azure CLI - To host chatbot

## 1.3. Data Source

Dataset is collected from - data-flair

url : https://drive.google.com/file/d/1763Y5zy7HmRYsOoBLQgUxQRGY6xCgQiN/view

(FIle used - intents.json)

## 1.4. Files to be used

- intents.json
- train_chatbot.py & chatgui.py
- Words.pkl & Classes.pkl
- chatbot_model.h5

# 2. Initial Processes

## 2.1. Import and Load data

Import Necessary Packages

```python
1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # In[1]:
5  import nltk
6
7  # In[2]:
8  from nltk.stem import WordNetLemmatizer
9
10 # In[3]:
11 lemmatizer = WordNetLemmatizer()
12
13 # In[4]:
14 import json
15
16 # In[5]:
17 import pickle
18
19 # In[6]:
20 import numpy as np
21
22 # In[7]:
23 from tensorflow import keras
24
25 # In[8]:
26 from keras.models import Sequential
27
28 # In[9]:
29 from keras.layers import Dense, Activation, Dropout
30
31 # In[10]:
32 import random
```

Parse intents.json files to python and read it

```python
35 # In[11]:
36 words = []
37
38 # In[12]:
39 classes = []
40
41 # In[13]:
42 documents = []
43
44 # In[14]:
45 ignore_words = ['?', '!']
46
47 # In[15]:
48 data_file = open('intents.json').read()
49
50 # In[16]:
51 intents = json.loads(data_file)
52
```

## 2.2. Pre-process Data

Tokenize the data

```
54  # In[17]:
55  for intent in intents['intents']:
56      for pattern in intent['patterns']:
57
58          w = nltk.word_tokenize(pattern)
59          words.extend(w)
60          documents.append((w, intent['tag']))
61
62          # add to classes list
63          if intent['tag'] not in classes:
64              classes.append(intent['tag'])
65
```

Lemmatize the words

```
68  words = [lemmatizer.lemmatize(w.lower()) for w in words if w not in ignore_words]
69
70
71  # In[19]:
72  words = sorted(list(set(words)))
73
74  # In[20]:
75  classes = sorted(list(set(classes)))
76
77  # In[21]:
78  print (len(documents), "documents")
79
80  # In[22]:
81  print (len(classes), "classes", classes)
82
83
84  # In[23]:
85  print (len(words), "unique lemmatized words", words)
86
87
88  # In[24]:
89  pickle.dump(words,open('words.pkl','wb'))
90
91
92  # In[25]:
93  pickle.dump(classes,open('classes.pkl','wb'))
94
```
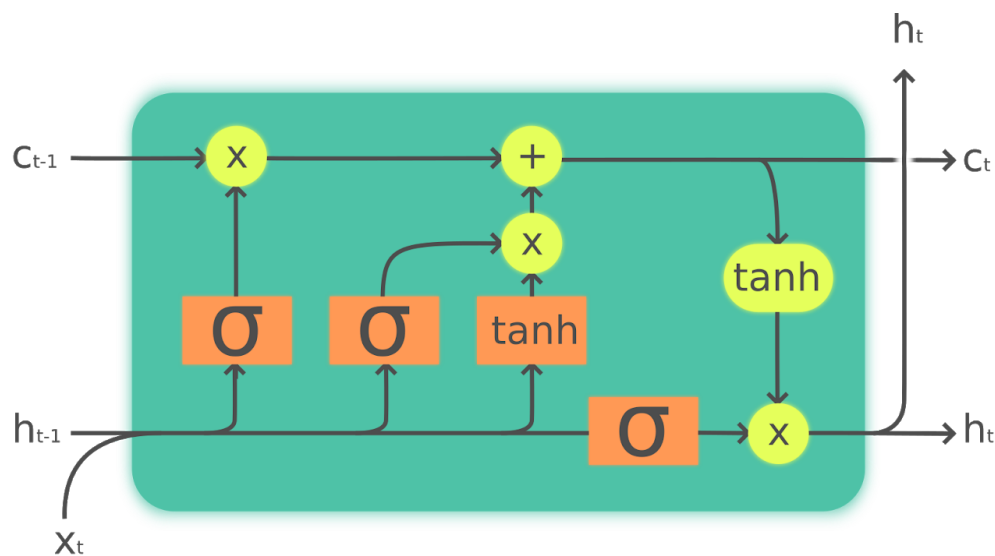
# 3. Train data and Build Model

## 3.1. About Algorithm

LSTM - Deep Learning model based on RNN architecture. Used to process sequences of data.

Commonly used in handwriting and speech recognition.

Mainly built to resolve long term dependency problem.

LSTM composes of 3 gates (input, output, forget) and a cell.



figure 1: Process of a cell in LSTM

Natural Language Understanding occurs via a bi-directional LSTM network.

An input is processed through three layers to bring out a conclusive result.

## 3.2. Training data and building the model

In training data, the input will be the pattern and output will be the class the pattern belongs to.

For the machine to understand patterns, text is converted to numbers

```
 99  #create training and testing datasets
100  training = []
101
102  # In[27]:
103  output_empty = [0] * len(classes)
104
105  # In[28]:
106  for doc in documents:
107      bag = []
108      pattern_words = doc[0]
109      pattern_words = [lemmatizer.lemmatize(word.lower()) for word in pattern_words]
110
111
112  # In[29]:
113  for w in words:
114      bag.append(1) if w in pattern_words else bag.append(0)
115
116      output_row = list(output_empty)
117      output_row[classes.index(doc[1])] = 1
118      training.append([bag, output_row])
119
120
121  # In[30]:
122  random.shuffle(training)
123
124  # In[31]:
125  training = np.array(training)
126
127  # In[32]:
128  train_x = list(training[:,0])
129
130
131  # In[33]:
132  train_y = list(training[:,1])
133
```

Predict the class of response using text-preprocessing

```
200 # In[59]:
201 def clean_up_sentence(sentence):
202     # tokenize the pattern - split words into array
203     sentence_words = nltk.word_tokenize(sentence)
204
205     # stem each word - create short form for word
206     sentence_words = [lemmatizer.lemmatize(word.lower()) for word in sentence_words]
207     return sentence_words
208
209 # In[60]:
210 def bow(sentence, words, show_details=True):
211     # tokenize the pattern
212     sentence_words = clean_up_sentence(sentence)
213
214     # bag of words - matrix of N words, vocabulary matrix
215     bag = [0]*len(words)
216     for s in sentence_words:
217         for i,w in enumerate(words):
218             if w == s:
219                 # assign 1 if current word is in the vocabulary position
220                 bag[i] = 1
221                 if show_details:
222                     print ("found in bag: %s" % w)
223     return(np.array(bag))
224
225 # In[61]:
226 def predict_class(sentence, model):
227     # filter out predictions below a threshold
228     p = bow(sentence, words,show_details=False)
229     res = model.predict(np.array([p]))[0]
230     ERROR_THRESHOLD = 0.25
231     results = [[i,r] for i,r in enumerate(res) if r>ERROR_THRESHOLD]
232
233     # sort by strength of probability
234     results.sort(key=lambda x: x[1], reverse=True)
235     return_list = []
236     for r in results:
237         return_list.append({"intent": classes[r[0]], "probability": str(r[1])})
238     return return_list
```

On predicting the class of response, random message is generated from relevant list of intents

```
241 # In[62]:
242 def getResponse(ints, intents_json):
243     tag = ints[0]['intent']
244     list_of_intents = intents_json['intents']
245     for i in list_of_intents:
246         if(i['tag']== tag):
247             result = random.choice(i['responses'])
248             break
249     return result
250
251 # In[63]:
252 def chatbot_response(text):
253     ints = predict_class(text, model)
254     res = getResponse(ints, intents)
255     return res
256
257 # In[64]:
258 import tkinter
```

Create GUI and execute the model

```
257 # In[64]:
258 import tkinter
259
260 # In[65]:
261 from tkinter import *
262
263 # In[66]:
264 def send():
265     msg = EntryBox.get("1.0",'end-1c').strip()
266     EntryBox.delete("0.0",END)
267     if msg != '':
268         ChatLog.config(state=NORMAL)
269         ChatLog.insert(END, "You: " + msg + '\n\n')
270         ChatLog.config(foreground="#442265", font=("Verdana", 12 ))
271         res = chatbot_response(msg)
272         ChatLog.insert(END, "Bot: " + res + '\n\n')
273         ChatLog.config(state=DISABLED)
274         ChatLog.yview(END)
275 base = Tk()
276 base.title("Hello")
277 base.geometry("400x500")
278 base.resizable(width=FALSE, height=FALSE)
279
280 # In[67]:
281 ChatLog = Text(base, bd=0, bg="white", height="8", width="50", font="Arial",)
282
283 # In[68]:
284 ChatLog.config(state=DISABLED)
285
286 # In[69]:
287 scrollbar = Scrollbar(base, command=ChatLog.yview, cursor="heart")
288 ChatLog['yscrollcommand'] = scrollbar.set
289
290 # In[70]:
291 SendButton = Button(base, font=("Verdana",12,'bold'), text="Send", width="12", height=5,
292                     bd=0, bg="#32de97", activebackground="#3c9d9b",fg='#ffffff',
293                     command= send )
294
295 # In[71]:
296 EntryBox = Text(base, bd=0, bg="white",width="29", height="5", font="Arial")
297
298 # In[72]:
299 scrollbar.place(x=376,y=6, height=386)
300 ChatLog.place(x=6,y=6, height=386, width=370)
301 EntryBox.place(x=128, y=401, height=90, width=265)
302 SendButton.place(x=6, y=401, height=90)
303
304 # In[73]:
305 base.mainloop()
306
307 # In[74]:
308 get_ipython().run_line_magic('run', 'train_chatbot.py')
309
310 # In[75]:
311 get_ipython().run_line_magic('run', 'chatgui.py')
```

# 4. Discussion

LSTM consists of memory cells which can store data for a certain period. Three gates control when data is input, output or forgotten. This helps LSTM learn longer term dependencies, neutralizing the vanishing gradient problem. This is why LSTM is ideal for the chatbot constructed.

This is a 'Retrieval based' chatbot which uses predefined classes of data to predict best responses.