

Human Behaviour Understanding Using Machine Learning and  
Psychological Methods

# *Depression Detection in Social Media*

*Nour Eldin Alaa Badr*  
*Wilmar Adrian Granados Cardenas*

December 19, 2020

## **Abstract**

Language plays a central role for psychologists who rely on manual coding of patient language for diagnosis. Using language technology for psychological diagnosis, particularly for depression detection could lead to low-priced screening test and affordable treatment. For Human Behaviour Understanding Using MachineLearning and Psychological Methods seminar's project we explore two neural architectures: Recurrent Neural Net-work (RNN) and Bidirectional Encoder Representations from Transformer (BERT), for depression detection on Reddit posts. In this report we will describe the data preparation and pre-processing as well as the implementation of the RNN and BERT models for detecting depression.

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Dataset Overview</b>	<b>5</b>
<b>3</b>	<b>Libraries used</b>	<b>5</b>
3.1	Why did we choose Pytorch instead of TensorFlow . . . . .	5
3.2	Why did we choose Spacy instead of NLTK . . . . .	6
<b>4</b>	<b>Folders and XML files parsing</b>	<b>6</b>
<b>5</b>	<b>Data pre-processing</b>	<b>8</b>
<b>6</b>	<b>Language Models</b>	<b>12</b>
6.1	RNN approach . . . . .	12
6.1.1	Introduction . . . . .	12
6.1.2	RNN Model . . . . .	13
6.1.3	Depression detection network architecture . . . . .	14
6.1.4	Training the network . . . . .	18
6.1.4.1	Loss and optimization functions . . . . .	19
6.1.4.2	Model hyperparameters . . . . .	19
6.1.4.3	Main steps of training the model . . . . .	20
6.1.4.4	Cross validation . . . . .	22
6.1.4.5	Results and observations . . . . .	23
6.1.5	Next steps regarding RNN . . . . .	25
6.2	Bert approach . . . . .	25
6.2.1	Introduction . . . . .	25
6.2.2	Bert Model . . . . .	26
6.2.3	Tokenization with Bert . . . . .	26
6.2.4	BertClassifier . . . . .	31
6.2.4.1	Bert Hyper-parameters . . . . .	31
6.2.4.2	Bert Training and Validation Functions . . . . .	32
6.2.4.3	Training and Results . . . . .	35
6.2.5	Insights gained . . . . .	36
	<b>References</b>	<b>39</b>

# 1 Introduction

Major Depressive Disorder (MDD), or commonly known as depression, is a health condition characterized by persistent sadness and lack of interest for previously enjoyable activities, feelings of guilt, disturbed sleep, tiredness, and poor concentration. These affections become depression symptoms when their occurrences last for at least 2 weeks [12]. Depending on the severity of the symptoms, psychologists classify depression as mild, moderate, or severe. The World Health Organization (WHO) considers depression a critical disease, as it might lead to suicide[20]. According to a Global Health Metrics article, more than 264 million people all over the world suffer from depression[7].

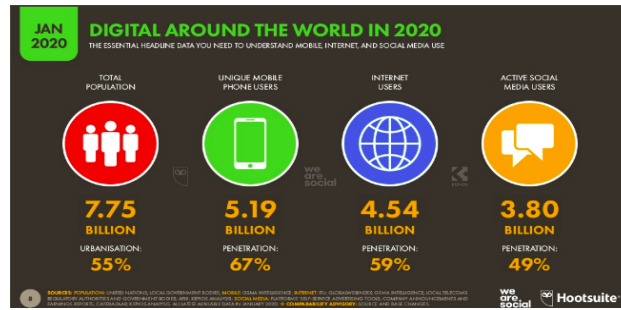


Figure 1: Global Reportal Slide (8) [8]

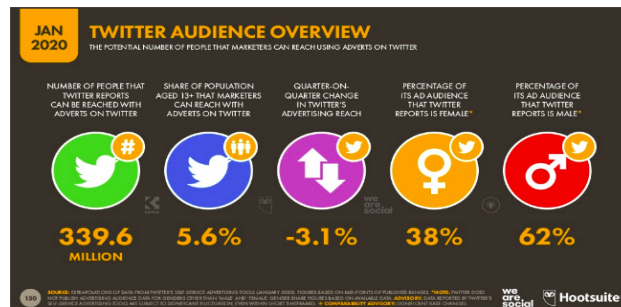


Figure 2: Number of twitter users slide (150) [8]

For clinical psychologists, language plays a central role in diagnosing depression and hence, many clinical instruments rely on manual coding of patient

language. Applying language technology in mental health diagnosis could lead to low-priced screening measures administered by more healthcare professionals. Researchers had begun targeting such issues using language technology to identify emotion in suicide notes, analyze the language of those with autistic spectrum disorders, and aid the diagnosis of dementia.[14] As depression increases on a global scale, it could become a global burden, which is why computer scientists and social researchers have cooperated to find ways for early detection of depression [2] and treatment [11] [5] using AI technologies. They found that AI can be a big help in detecting depression through various approaches, as it learns fast, can detect depression with high accuracy, can detect it from facial expressions, voices, and texts. As such, research on early detection of depression has shown relevant insight and brought us to the idea of detecting it in texts on social media.

## **2 Dataset Overview**

We have used a dataset from the Early Risk Prediction on the Internet (erisk 2017) workshop. This dataset provides user-generated content consisting of Reddit posts organized and processed chronologically. Every user has received a label as risk or non-risk (of depression). The dataset has two parts: a training dataset and a test dataset, and each part has ten chunks with series of XML files. These XML files store users' posts and the respective posts' comments.

## **3 Libraries used**

When it comes to training model in deep learning, our minds goes to two powerful tools which are Pytorch and Tensorflow. While in tokenization there are Spacy and NLTK. In this project we have used Pytorch library for training and we have used Spacy for tokenization.

### **3.1 Why did we choose Pytorch instead of TensorFlow**

Tensorflow has some advantage if it is choosen for training your model, such as it has a useful tool called Tensorboard which allows you to visualize your machine learning model, this tool isn't available in Pytorch. Also Tensorflow has a larger online community than Pytorch. Despite of all these advantages

we have decided to use Pytorch as it is more pythonic. By this we mean that you don't have to spend a lot of time learning how to use the library as the case in Tensorflow, you only need to know how to program in python which make it very easy and straight forward to use.

### 3.2 Why did we choose Spacy instead of NLTK

Our selection of SpaCy over NLTK has one main reason: performance. NLTK is a well-developed library for string processing. It takes a string as input and returns a string or a list (array) of strings. However, NLTK's array focus makes it a slower option than SpaCy for processing large amounts of text. SpaCy has an object-orient approach, meaning that it takes a string as input and returns a document object whose chunks and tokens are also objects. Spacy's OOP approach complements its industry-driven development that promotes its implementation with the fastest algorithms for text processing.

## 4 Folders and XML files parsing

This section focuses on folders and XML parsing. As shown in the following class diagram figure3, this class consists of 2 main functions:

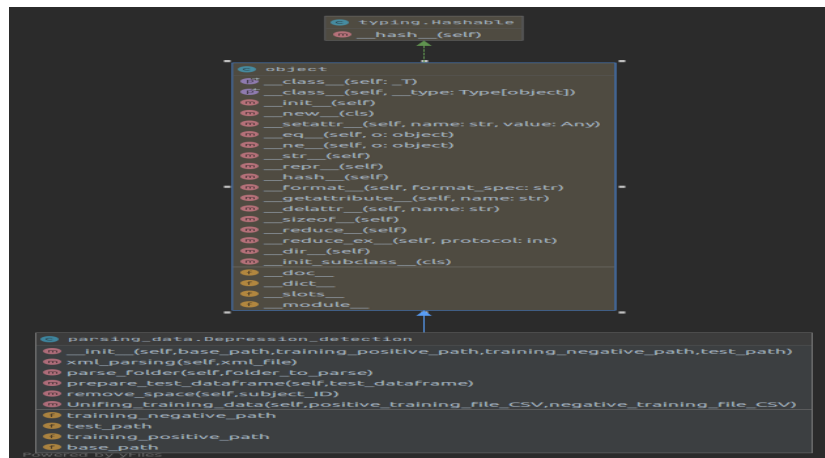


Figure 3: Folders and XML files parsing

- **Parsing Folder function:** is the function that parses the training/testing directory and the folders inside it in order to call the function that

parses the XML files. This is done by iterating on each directory and till finding a file and then call the XML parsing function

```

1 def parse_folder(self, folder_to_parse):
2     '''
3     1. It parses the training /testing directory where
4     each folder contains 10 chunks as well as
5     2. It calls the XML parser and returns list of data
6     frames
7
8     :param folder_to_parse:
9     :return: list of frames
10    '''
11    frames = []
12    folders_path = Path(folder_to_parse)
13    for directory in folders_path.iterdir():
14        if directory.is_dir():
15            print(directory.name)
16            for file in directory.iterdir():
17                if file.is_file() and not file.name.
18                startswith('.') and not file.name.startswith('.'):
19                    frame = self.xml_parsing(file)
20                    frames.append(frame)
21
22    return frames

```

**Listing 1:** Parsing Folders

- **Parsing XML function:** This function parses the xml files into a dataframe with 5 columns ID, TITLE(which is the post), DATE, INFO and the TEXT(which are the comments).

```

1 def xml_parsing(self, xml_file):
2     ''' Parses a xml file with the following structure
3     <INDIVIDUAL>
4     <ID>...</ID>
5     <WRITING>
6     <TITLE>...</TITLE><DATE>...</DATE><INFO>...</
7     INFO><TEXT>...</TEXT>
8     </WRITING>
9     ....
10    </INDIVIDUAL>
11    Returns DataFrame with columns: ID TITLE DATE INFO
12    TEXT
13    '''

```

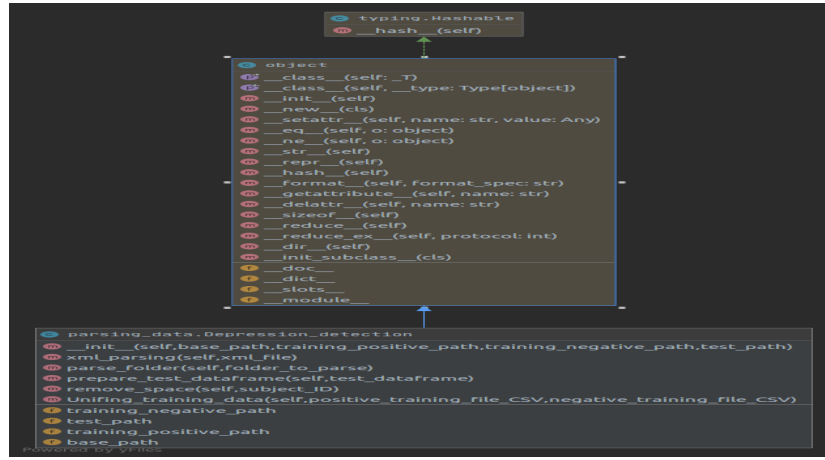


Figure 4: Data Pre-processing Class Diagram

```

12 xtree = et.parse(xml_file)
13 xroot = xtree.getroot()
14
15 subject_id = xroot.find('ID').text
16 writings = []
17 for writing in xroot.findall('WRITING'):
18     xml_data = {}
19     xml_data['ID'] = subject_id
20     xml_data['TITLE'] = writing.find('TITLE').text
21     xml_data['DATE'] = writing.find('DATE').text
22     xml_data['INFO'] = writing.find('INFO').text
23     xml_data['TEXT'] = writing.find('TEXT').text
24     writings.append(xml_data)
25 df = pd.DataFrame(writings)
26 return df
  
```

Listing 2: Parsing XML files

## 5 Data pre-processing

In this section, data will be preprocessed to feed it to our model. As shown in the following class diagram fig 4, there are 4 main functions:

- **Tokenization:** This function removes https links and punctuation, removes spaces and numbers, lowering the case, tokenize the whole



text and finally save them in a pickle file to use them later. We had two options to tokenize our posts. The first option is to tokenize it without using any libraries at all, just using splitting function. The second option is to use libraries which would be more efficient as these libraries are more efficient and trained to do so.

```
1
2 def tokenization(self, unified_training_df, train=False):
3     '''
4     Preprocess the whole text ie.
5     * remove https links and punctuation
6     * remove spaces and numbers
7     * lowering the case and finally tokenize it.
8     :param unified_training_df:
9     :return: tokens
10    '''
11    #get punctuations from string
12    punctuations = string.punctuation
13    corpus = unified_training_df.TITLE_TEXT
14    corpus = corpus.to_list()
15
16    #Take each line of the subject and combine all of
17    #them in one big text(corpus)
18    text = []
19    for x in corpus:
20        text.append(x)
21
22    text_joined = ' '
23    text = text_joined.join(text)
24
25    #remove https links and punctuations
26    text = re.sub(r'http\S+', '', text)
27    all_text = ''.join([c for c in text if c not in
28    punctuations])
29
30    #Tokenize the whole text using spacy tokenizer
31    words = self.tokenizer(all_text)
32
33    #delete space and numbers and lower case all the
34    #tokens
35    tokens = [token.lower_ for token in words if not
36    token.is_space and not token.like_num]
```

```
36     return tokens
```

### Listing 3: Tokenization

- **Downsampling:** This function focuses on solving the imbalanced data problem, which is that our dataset has more positive training data than negative data. We had two options: either an oversampling technique by repeating negative data couple of times to equalise its size with the positive data or using a downsampling technique, in which we delete data from positive data. Oversampling will cause over-fitting and might decrease the accuracy due to repeated data. On the other hand, downsampling will lead to losing data and information from our dataset. To balance the data, we chose the technique of downsampling, because oversampling will cause over-fitting.

```
1 def downsampling(self, unified_training_df_preprocessed):
2     '''
3     downsampling majority(un)
4     :param unified_training_df_preprocessed:
5     :return:
6     '''
7     # separate minority and majority classes
8     un_depressed = unified_training_df_preprocessed[
unified_training_df_preprocessed.LABEL == 0]
9     depressed = unified_training_df_preprocessed[
unified_training_df_preprocessed.LABEL == 1]
10
11     undepressed_downsampled = resample(un_depressed,
12                                       replace=False, #
13                                       sample without replacement
14                                       depressed), # match minority n
15                                       n_samples=len(
16                                       depressed),
17                                       random_state=27)
18     # reproducible results
19
20     # combine minority and downsampled majority
21     downsampled = pd.concat([undepressed_downsampled,
22                              depressed])
23
24     # checking counts
25     downsampled.LABEL.value_counts()
26
27     #save it to csv file
```

```
23 downsampler.to_csv('./downsampled_data.csv')
```

**Listing 4:** Downsampling

- **Encoding the words:** This function gives each token a unique integer. It is done by mapping each word to an integer using dictionaries. We are taking this step as embedding lookup in Pytorch requires an integer to the network. Embedding lookup or lookup table is embedding weight matrix, which consists of vectors for each word. This part will be discussed more in detail in the embedding layer section.

```
1 def vocab_to_int(self, tokens):
2
3     '''
4     Annotate all the vocabs in the text with a
5     corresponding integer.
6     :return: self.vocab_to_int
7     '''
8
9     ## Build a dictionary that maps words to integers
10    counts = Counter(tokens)
11    vocab = sorted(counts, key=counts.get, reverse=True)
12    # print(vocab[:30])
13    # 1 in enumerate means the dictionary will begin with
14    # one instead of 0
15    self.vocab_to_int = {word: ii for ii, word in
16                        enumerate(vocab, 1)}
17
18    return self.vocab_to_int
```

**Listing 5:** Vocab\_to\_integer

- **Padding sequence:** To train our model, we need all the inputs of the neural network to have the same shape and size, but in our dataset we have found that the posts and the comments of the subjects don't have the same length. This is why we had to create a function to deal with both short and very long posts or comments to make all of them have the same length.

If we have some short posts, the function will pad them with zeros to create the same size of the sequence length we decide to use. Otherwise, if the posts of the subjects are too long the function will truncate

them to also have the same length as the sequence length.

```
1 def pad_features(self, text_integers, seq_length):
2     ''' Return features of text_ints, where each text is
3     padded with 0's
4     or truncated to the input seq_length.
5     '''
6     ## implement function
7     # getting the correct rows x cols shape
8     features = np.zeros((len(text_integers), seq_length),
9     dtype=int)
10    # for each review, I grab that review
11    for i, row in enumerate(text_integers):
12        # this will take row and last columns and leave
13        # first columns with zeros (check last cell to
14        # understand) , then [:seq_length] will shrink long data
15        features[i, -len(row):] = np.array(row)[:
16        seq_length]
17    #convert features to numpy
18    features = np.array(features)
19    return features
```

**Listing 6:** Padding\_features

After Encoding and padding features, we are now ready for the next step, which is feeding these features to our model.

## 6 Language Models

This section presents the two models implemented to detect depression from online posts. We will give a description of each implementation with examples of the text's representation used in each model.

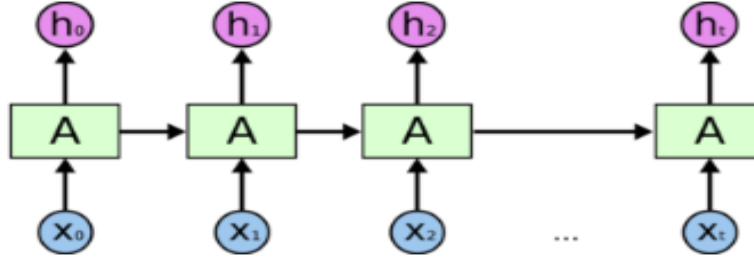
### 6.1 RNN approach

#### 6.1.1 Introduction

This section will give a general description on RNN(Recurrent Neural Network) Model, how we will implement the RNN model in Pytorch and how to use it in detecting depression.

### 6.1.2 RNN Model

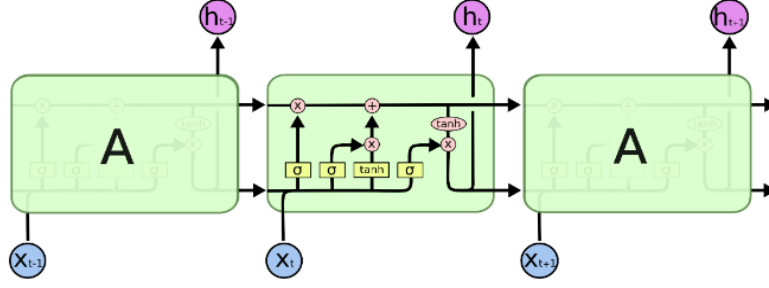
As the main purpose of this report is to explain how we implemented RNN rather than explaining RNN itself, we will only give a general explanation about it. We can simply describe RNN as a normal feedforward neural network, but it only differs in that it has internal memory. As shown in figure 5, RNN has two inputs, it has the current input and the previous input, which it gets from the output of the previous cell. Usually the output is referred to as hidden state [15].



**Figure 5:** recurrent neural network [17]

In our project, we have decided to use LSTM(Long Short Term Memory), which is an extension of the normal(vanilla) RNN. We have found that RNN has a problem called vanishing gradient. It has been noticed that during backpropagation in RNN, the gradient (which is the update of neural network weights) shrinks through time. In this case the gradient value becomes too small and it doesn't contribute much to the learning, especially recurrent neural networks that get small gradient update [6]. That is why we have decided to use LSTM.

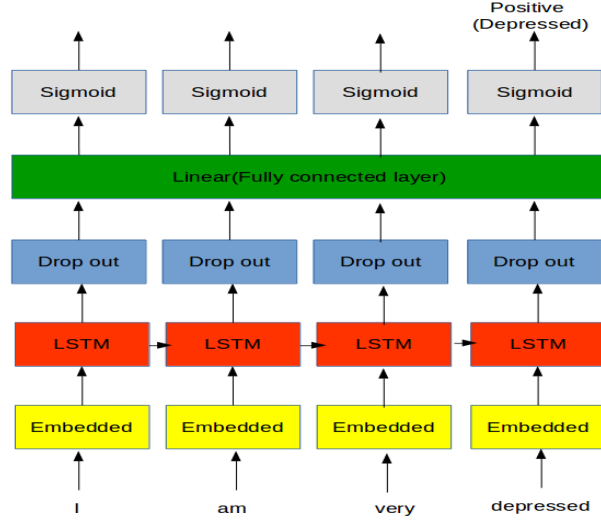
LSTM as shown in figure 6, briefly has two outputs instead of just one. One output is the cell state, which is the memory of the LSTM that has all the previous information, the other output is the hidden state which is the output of the cell [6].



**Figure 6:** Long Short Term Memory [17]

### 6.1.3 Depression detection network architecture

As shown in our figure 7, our network architecture consists of 5 layers: Embedding layer, LSTM layer, Dropout layer, Linear layer and Sigmoid layer. We will discuss each layer in detail:



**Figure 7:** Depression detection network architecture

## 1. Embedding Layer:

To train our model and pass the inputs to our RNN (LSTM), our features should be represented in vectors, as neural networks need a vector as an input. So we have thought of using one hot encoded vectors. However, one hot vector encoder has two disadvantages:

- It contains ones and zeros only, which makes it difficult for the neural network to learn semantic representations.
- Usually, if you have a large number of words in your vocab, it might not be a good idea to use one hot vector encoder, since you will have a matrix with high dimensionality, which will make the model train slowly.

As we have seen these disadvantages, we have decided to use embedding layer to increase efficiency, reduce dimensionality and each vector has meaningful values(weights) to the model [10]

To make it more clear, let us have a simple example. Let us assume that we have two posts from users:

*My day was very depressing*

*Today was fine*

firstly we will tokenize these two sentences to be as following:

```
1 [My, day, was, very, depressing]
2 [Today, was, fine]
```

**Listing 7:** Tokenization example

Next step is to give each word a unique number as well as applying padding to sentences(posts), assuming that sequence length = 5. We have to note that **was** token is repeated in both sentences(posts) so they will have the same unique integer.

```
1 [1, 2, 3, 4, 5]
2 [0, 0, 0, 6, 3, 7]
```

**Listing 8:** Encoding and Padding example

Now we are ready to create our embedding weight matrix. Pytorch has made it easy for us to create these embedding weight matrix as it creates it and can fill and optimize these embedding vectors during

training phase. Alternatively, we can fill them from a pretrained model like Glove.

The expected size of this lookup table(embedding weight matrix) is (number of unique words+ 1 , embedding dimension), where the 1 added to the unique words is due to the zero padding. So in our example it will be (8,3), if we assume that the embedding dimension is 3.

```

1 [in] : self.embedding = nn.Embedding(vocab_size,
    embedding_dim) #(8,3)
2
3
4 [out] : #Lookup table
5 index      Embedding
6 0          [2.5, 1.7, 1.9]
7 1          [2.9, 1.3, 1.1]
8 2          [1.9, 2.3, 2.1]
9 3          [1.1, 2.9, 2.5]
10 4         [1.8, 2.2, 2.8]
11 5         [3.8, 1.2, 3.2]
12 6         [1.8, 2.2, 2.4]
13 7         [4.2, 3.7, 2.4]

```

**Listing 9:** Embedding Layer

As shown in the lookup table, each word now has a unique index as well as an embedded vector.

## 2. LSTM layer:

In this section, we will focus more on the programmatic part than the theoretical part, as we have discussed this part theoretically in the previous section.

To create the LSTM layer, we have to give the function the following arguments. It takes the input size, the hidden dim, the number of lstm layers, dropout probability and setting batch first to true to get the following output dimensions (batch\_size,sequence\_length,input dimension). After creating the layer, we feed it with inputs, which are the embeddings and the hidden state.

```

1 self.lstm = nn.LSTM(embedding_dim, hidden_dim, n_layers,
    dropout=drop_prob, batch_first=True)
2

```



```
3 lstm_out, hidden = self.lstm(embeds, hidden)
```

**Listing 10:** LSTM Layer

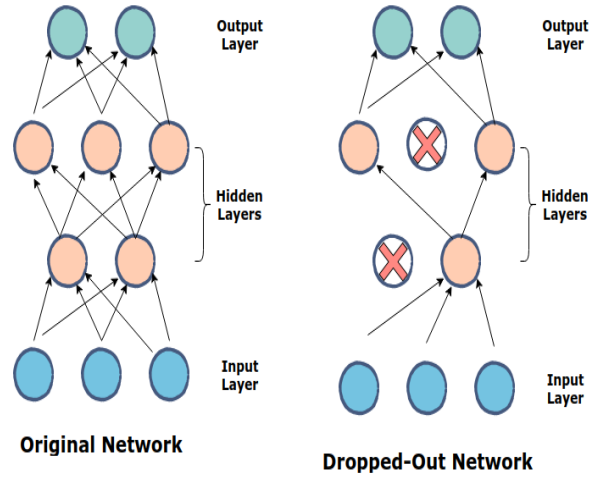
As we discussed earlier, we have hidden and cell states in LSTM and we need to initialize their weights at the beginning of each epoch or beginning of training. During our implementation we have tried two approaches. Firstly, we have tried to initialise the weights with zeros, then we tried to initialise the weights randomly. We have found that the accuracy has increased by 1% to 2% .

```
1 def init_hidden(self, batch_size):
2     ''' Initializes hidden state '''
3     # Create two new tensors with sizes n_layers x
4     batch_size x hidden_dim,
5     # initialized to zero, for hidden state and cell
6     state of LSTM
7     weight = next(self.parameters()).data
8
9     if (self.train_on_gpu):
10         device = torch.device("cuda:0") # Uncomment this
11         to run on GPU
12         hidden = (torch.randn(self.n_layers, batch_size,
13                                self.hidden_dim, device=device),
14                   torch.randn(self.n_layers, batch_size,
15                                self.hidden_dim, device=device))
16     return hidden
```

**Listing 11:** hidden weights

### 3. Dropout layer:

We are using this layer in order to prevent overfitting. This happens because, while we are updating our weights during training, they become dependent on the dataset that we are using [16]. So, during training we dropout a certain value from our network as shown in figure 8. We have tried to dropout 30 percent and 50 percent; there was not much difference in accuracy between them.



**Figure 8:** Dropout layer [19]

#### 4. Fully connected layer(Linear layer):

This is used to map our LSTM layer output to the output size we need. This means that it takes the hidden dimensions that we have chosen as an input and output only one value.

```
1 (fc): Linear(in_features=256, out_features=1, bias=True)
```

**Listing 12:** fully connected layer(Linear layer)

#### 5. Sigmoid layer:

This will make all the output values we have either 0 or 1. As shown in figure 7, we are only interested in the last sigmoid output on the right as we are not expecting a sequence output from the other sigmoid outputs.

**\*\*NOTE:\*\*** Please check our github link for the full network architecture code.

### 6.1.4 Training the network

Before starting the training process, we have to choose our loss and optimization functions and adjust our hyparameters.

#### 6.1.4.1 Loss and optimization functions

##### 1. Loss function:

Loss functions is telling us how wrong our predictions are during our training process. In other words, we are using it to evaluate how well our algorithm models our data. If the loss is high, then we have bad accuracy or bad performance. In our project we are using **Binary Cross Entropy Loss(BCE loss)**. We are using BCE loss as we have a single output which is 0 or 1 as well as BCE loss in pytorch was designed to work with sigmoid outputs.

```
1 criterion = nn.BCELoss()
```

**Listing 13:** Loss function

##### 2. Optimizer function:

This function simply defines how neural network learn and by that we mean that it finds the values of parameters, such that the loss function is at its lowest value.

In our project we have used Adam optimization, which is just an extension from Stochastic gradient decent. We have used it as one of the newest optimizers, as it has many benefits according to their authors [9]. It has efficient computations, requires little memory and works well for problems with very high noise gradients.

```
1 optimizer = torch.optim.Adam(RNN_net.parameters(), lr=lr)
   #lr is learning rate
```

**Listing 14:** Adam optimizer

#### 6.1.4.2 Model hyperparameters

1. **Vocab size** : It defines the size of our input, which is the word tokens.
2. **Output size** : It is the number of outputs in our model, which is 1 in our case as we have one output, either zero or one.

3. **Embedding dimension** : It is the size of our embeddings, which is the number of columns in an embedding vector which we have explained earlier. It is also known as the number of columns in embedding lookup table. It usually has a value of 200, 300 or 400.
4. **hidden dimension** : In our LSTM cells, we have a certain number of units in the hidden layers. These values could be 128, 256 or 512. Sometimes the more the dimension increases, the more time it takes to train the model, the better performance we get.
5. **number of layers** : It is the number of the LSTM layers in a network. It is usually between 1 and 3 layers.
6. **learning rate** : It is the step size at each iteration and, as shown in the optimization section, learning rate is a tuning parameter in the optimization algorithm.
7. **epochs** : It is the number of iterations through the whole training dataset.

```

1 # Instantiate the model w/ hyperparams
2 vocab_size = len(vocab_to_ints_training)+1 # +1 for the 0
   padding + our word tokens
3 output_size = 1
4 embedding_dim = 200
5 hidden_dim = 256 #128
6 n_layers = 1 #2
7 lr = 0.0001 #0.01
8 epochs = 5 #10
9

```

**Listing 15:** Model hyperparameters

#### 6.1.4.3 Main steps of training the model

1. At the beginning of each epoch during training, we will call the function that initializes the hidden state.

```

1 h = RNN_net.init_hidden(batch_size)

```

**Listing 16:** hidden state initialization

2. At the beginning of the batch loop, we will try to clear all the old gradients we have from the last step. If we didn't do this step, all the gradients would be accumulated from the step `loss.backward()` call, which we will discuss in further steps. Hence, in this step we will zero the accumulated gradients.

```
1 # zero accumulated gradients
2 RNN_net.zero_grad()
```

3. We will perform feed forward propagation by passing the inputs to the model to get predicted outputs.

```
1 # get the output from the model
2 output, h = RNN_net(inputs, h)
```

**Listing 17:** Feedforward propagation

4. In this step, we will calculate the loss between the predicted output and the actual output. After calculating the loss we will backpropagate to compute the derivatives of the loss.

```
1 # calculate the loss and perform backprop
2 loss = criterion(output.squeeze(), labels.float())
3 loss.backward()
```

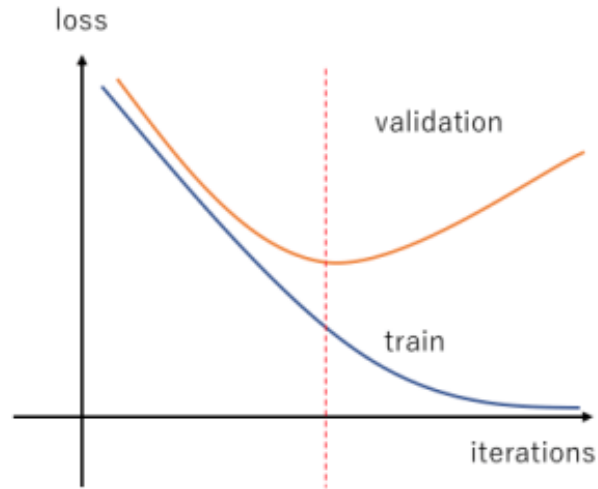
**Listing 18:** Loss calculation and perform backpropagation

5. In this step, we will update our parameters by performing a single optimization step.

```
1 optimizer.step()
```

**Listing 19:** Optimizer step

6. Lastly, our model will be saved under one condition: the validation loss is decreasing. If the past validation loss value is more than the newest validation value, we will not save the model. This last step helps us to stop our model from overfitting, as during training the training loss decreases aside from the validation loss, but at one breaking point the validation loss starts to increase again as shown in figure 9.



**Figure 9:** Train loss vs validation loss [13]

```

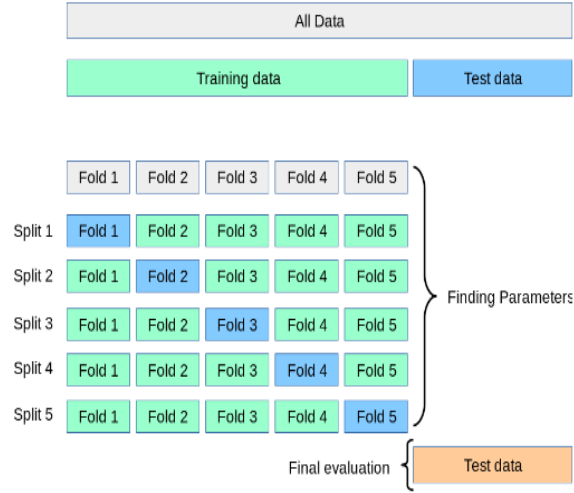
1 if np.mean(val_losses) <= valid_loss_min:
2     print('Validation loss decreased ({:.6f} --> {:.6f}).
3       Saving model ...'.format(
4         valid_loss_min,
5         np.mean(val_losses)))
6     torch.save(RNN_net.state_dict(), '
model_trained_RNN_not_pretrained.pt')
    valid_loss_min = np.mean(val_losses)

```

**Listing 20:** Saving Model

#### 6.1.4.4 Cross validation

In order to get a good evaluation of the model and evaluate its performance, we need to split our training data into folds as shown in figure 10. Each time we split our data, there should be a fold that is not fed into the training dataset, but fed into our validation dataset to evaluate our model. We thus try to train the model each time on different folds to see how accurate it is while working on unseen data, as well as trying to avoid overfitting or underfitting.



**Figure 10:** cross validation [3]

#### 6.1.4.5 Results and observations

As we have discussed earlier, we have trained our model using cross validation using different splits for training and validation. As shown in the following table 1, we got an accuracy with around 74 % during training, but when we have tried to test our model on unseen data we get an accuracy of 64.4 %. As shown in the table 1 the second split has higher accuracy on unseen data. The model is trained on the following hyperparameters:

```

1 output_size = 1
2 embedding_dim = 200
3 hidden_dim = 256
4 n_layers = 1
5 lr = 0.0001
6 epochs = 5

```

**Listing 21:** CV Hyperparameters

Epochs	Cross Validation									
	Split 1		Split 2		Split 3		Split 4		Split 5	
	Accuracy during training	Accuracy on unseen data(test data)	Accuracy during training	Accuracy on unseen data(test data)	Accuracy during training	Accuracy on unseen data(test data)	Accuracy during training	Accuracy on unseen data(test data)	Accuracy during training	Accuracy on unseen data(test data)
1	54%	60.8%	55%	64.4%	57%	60%	56%	61%	53%	59.4%
2	63.9%		65%		63%		62.3%		63.2%	
3	65.1%		74.6%		72.1%		73%		64.5%	
4	66%		74.7%		72.3%		73.2%		66.2%	
5	66.9%		74.7%		72.4%		73.2%		66.3%	

**Table 1:** Data trained on different splits (Cross Validation)

While training the model, we have first tried to train the model on 10 epochs. What we noticed is that both training loss and validation loss decrease till epoch 5. After that, the model starts to overfit and the validation loss starts to increase again. So we have decided to train the model with 5 epochs only in order to avoid overfitting. Please refer to figure 9.

We have tried to increase the embedding layers from 200 to 300 as well as increasing the hidden dimensions from 256 to 512 besides increasing the number of layers from 1 to 2. We have noticed that as we increase the hidden dimensions and the number of layers, the time the model takes for training almost doubled. The accuracy didn't change that much, we have reached an accuracy of 67% during training and around 59% on unseen test data. We have tried also to increase the number of layers to 3 and the embedding dimensions to 400, we got an accuracy of 75% during training and 62% on unseen test data. A summary of our training can be found in the following table 2.

	Embedding dimensions = 200 , Hidden dimension = 256, Number of layer = 1, Epoch = 0.001	Embedding dimensions = 300 , Hidden dimension = 512, Number of layer = 2, Epoch = 0.001	Embedding dimensions = 400 , Hidden dimension = 512, Number of layer = 3, Epoch = 0.001
During training accuracy	74%	67%	75%
Testing on unseen data accuracy	64.4%	59%	62%

**Table 2:** (Accuracy Comparison)



We have also tried to train our network using pretrained weights from a pretrained model. We have used Glove pretrained word vectors, which offers pretrained weights based on tokens from twitter. To set the model to train from the pretrained model, we first have to only look for the tokens that we have in our vocab within Glove tokens and extract the weights of our tokens only. Then we have to call the pretrained method from pytorch during our training.

```
1 elif pretrained == True:
2     num_embeddings, embedding_dim = weights_matrix.shape
3     self.embedding = nn.Embedding(num_embeddings,
4     embedding_dim)
5     self.embedding.load_state_dict({'weight': weights_matrix
6     })
```

**Listing 22:** Pretrained model

We have trained the model with pretrained weights using the same hyperparameters in listing 21 .Surprisingly using pretrained model didn't change much in the accuracy, as we got around 70 % accuracy during training and around 60% during testing it on unseen data.

### 6.1.5 Next steps regarding RNN

We would be interested to enhance our model accuracy using RNN, we will try different hyperparameters, maybe we can try to replace LSTM with GRU instead . We have tried also to use Skorch which is a library that allows us to use Pytorch with Sklearn . One of the advantages of this library that during training it applies different hyperparameters randomly as well as it train on GRU and LSTM both together. But unfortunately it gives shows the highest accuracy without stating which hyperparameters gave this accuracy. We would like to invest more time on using this library as it might help a lot. We got an accuracy of 69% using Skorch.

## 6.2 Bert approach

### 6.2.1 Introduction

This section will give a general description of the Bidirectional Encoder Representations from Transformers (BERT) language model's operation, its main

components, and its implementation for detection of depression. This description will indicate how a better understanding of the language model functioning could contribute to enhanced implementations.

### 6.2.2 Bert Model

BERT is a language model developed by Google and designed to carry out different tasks. The base model was trained on Wikipedia and the Google Books Corpora; both data sets sum up to 3500 million words. This amount of data gives BERT the flexibility for a broad range of applications. BERT requires fine-tuning to perform specific tasks, such as question answering, natural language understanding, sentiment analysis, and depression detection.

Transformers are BERT's base architecture; they replace recurring layers (LSTM) with attention layers [18]. Attention layers encode each sequence's token based on its contextual relation with other sequence's tokens. This encoding introduces the context of each token in the numerical representation of whole sequence. For this type of encoding, the transformers models are also called contextual embeddings. Transformers models also generate positional embedding, which, like the contextual embedding, allows knowing each token's relative position in the encoded text sequence.

Transformer-based models' implementation consists of two phases. The first phase is called pre-training [1]. In this phase, the model learns language structures and generic knowledge of each word's meaning. The learning occurs through "exercises", in which the model has to predict a word or words omitted in a sentence. In our project, we skipped the pre-training phase using BERT, which has been pre-trained in two large databases, as it has been mentioned.

Next, we will explain the general operation of BERT:

### 6.2.3 Tokenization with Bert

```
1 PRE_TRAINED_MODEL_NAME="bert-base-cased"
```

```

2 tokenizer=BertTokenizer.from_pretrained(
    PRE_TRAINED_MODEL_NAME)

```

### Listing 23: Pre-trained Model

BERT base-cased model requires its input as a sequence divided into two parts. The first part of the sequence will have the original text, and the second part will have pad tokens that indicate it is empty. After BERT's pre-training, the first token of the sequence condenses the whole sequence representation, for the tokenizer class uses the pre-trained model, which has 100 million parameters and differentiates between upper and lower case written words. The following example will show the tokenization process in a text sample:

```

1 txt_sample = """I am sick of the Hall and the hill, I am sick
    of the moor and the main. Why should I stay? can a
    sweeter chance ever come to me here? O, having the nerves
    of motion as well as the nerves of pain, Were it not wise
    if I fled from the place and the pit and the fear?"""
2
3 # The Valley, (Part One, I, xvi) (Talonbooks, 2014), by Joan
    MacLeod.

```

### Listing 24: Text Sample

Since the Reddit posts have a confidentiality agreement to respect users' privacy, a fragment of the novel the Valley by Joan MacLeod has been taken as an example.

By using the tokenizer in the sample text, the following tokens are obtained:

```

1 txt_token = tokenizer.tokenize(txt_sample)
2 print('Tokens: ', txt_token)
3
4 Tokens:  ['I', 'am', 'sick', 'of', 'the', 'Hall', 'and', 'the',
    ', ', 'hill', ', ', 'I', 'am', 'sick', 'of', 'the', 'm', '##',
    'oor', 'and', 'the', 'main', '.', 'Why', 'should', 'I', ' ',
    'stay', '?', 'can', 'a', 'sweet', '##er', 'chance', 'ever',
    'come', 'to', 'me', 'here', '?', 'O', ', ', 'having', 'the',
    ', ', 'nerves', 'of', 'motion', 'as', 'well', 'as', 'the', ' ',
    'nerves', 'of', 'pain', ', ', 'Were', 'it', 'not', 'wise', ' ',
    'if', 'I', 'fled', 'from', 'the', 'place', 'and', 'the', ' ',
    'pit', 'and', 'the', 'fear', '?']

```

### Listing 25: Tokens

Tokenization shows difficulty for a psychological interpretation of depression detection based on the documents' representation. As it can be seen, in tokenized text, the text's semantic structure is altered by including masked tokens like 'm', '##oor', '##er'. These tokens give advantages for the text's processing, since their representation allows generalizing and better performance, as explained in BERT's launching paper. "The masked language model [MLM] randomly masks some of the tokens from the input, and the objective is to predict the original vocabulary id of the masked word based only on its context. Unlike left-to-right language model pre-training, the MLM objective allows the representation to fuse the left and the right context, which allows us to pre-train a deep bidirectional Transformer." [4]. The tokenization benefits the implementation of the transformers' architecture, but shuns meaningful psychological research clues. It should be investigated if the original text's semantic structure changes follow any regularity according to the type of tokens that are altered. Transformers language models give promising results, but these come at the cost of limited or none psychological interpretability.

After the tokenization, each token received a numerical assignment. The numeration comes from the pre-trained model, and is necessary because neural networks cannot handle textual inputs.

```

1 txt_token_ids = tokenizer.convert_tokens_to_ids(txt_token)
2 print('Token IDS: ', txt_token_ids)
3
4 Token IDS: [146, 1821, 4809, 1104, 1103, 1944, 1105, 1103,
             4665, 117, 146, 1821, 4809, 1104, 1103, 182, 15626, 1105,
             1103, 1514, 119, 2009, 1431, 146, 2215, 136, 1169, 170,
             4105, 1200, 2640, 1518, 1435, 1106, 1143, 1303, 136, 152,
             117, 1515, 1103, 10846, 1104, 4018, 1112, 1218, 1112,
             1103, 10846, 1104, 2489, 117, 8640, 1122, 1136, 10228,
             1191, 146, 6192, 1121, 1103, 1282, 1105, 1103, 7172, 1105,
             1103, 2945, 136]
```

**Listing 26:** Tokens' numerical representation

The TokenizerClass, imported with BERT pre-trained model, has a function to encode the text to be processed by the model's transformers networks. The following function shows an implementation of encoding for the model re-training. This function takes as parameters: `txt_sample`, corresponding to the sample text seen before. In the model's implementation for depression detection, this parameter will take users' Reddit posts. `max_length` takes a

value of 100 tokens, in this example and the implementation for detecting depression. The following five boolean parameters: activate the truncation of the sequences, add of special tokens, leave the identifying token types, activate the padding, and returns the attention mask. The last parameter return the tensors.

```

1 encoding = tokenizer.encode_plus(
2     txt_sample,
3     max_length = 100,
4     truncation = True,
5     add_special_tokens = True,
6     return_token_type_ids = False,
7     pad_to_max_length = True,
8     return_attention_mask= True,
9     return_tensors= 'pt')

```

**Listing 27:** Encoding Functions

Checking the encoded text, we have the representation that BERT requires: a sequence with a special token at the first position '[CLS]' followed by the texts' token, a delimiting token '[SEP]' that closes the first part of the sequence. Next to this token, we have the second part of the sequence that will receive the padding token '[PAD]'. After the re-training, the '[CLS]' token will store the entire document's representation for classification.

```

1 print(tokenizer.convert_ids_to_tokens(encoding['input_ids']
2     [0]))
3 ['[CLS]', 'I', 'am', 'sick', 'of', 'the', 'Hall', 'and', 'the',
4  ', ', 'hill', ', ', 'I', 'am', 'sick', 'of', 'the', 'm', '##',
5  'oor', 'and', 'the', 'main', '.', 'Why', 'should', 'I', 'stay',
6  '?', 'can', 'a', 'sweet', '##er', 'chance', 'ever', 'come',
7  'to', 'me', 'here', '?', '0', ', ', 'having', 'the', 'nerves',
8  'of', 'motion', 'as', 'well', 'as', 'the', 'nerves', 'of',
9  'pain', ', ', 'Were', 'it', 'not', 'wise', 'if', 'I', 'fled',
10 'from', 'the', 'place', 'and', 'the', 'pit', 'and', 'the',
11 'fear', '?', '[SEP]', '[PAD]', '[PAD]', '[PAD]', '[PAD]',
12 '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]',
13 '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]',
14 '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]', '[PAD]',
15 '[PAD]']

```

**Listing 28:** Encoded Text Sample

In the tokens’ numerical assignment, we see that the tokenizer assigns to the tokens '[CLS]', '[SEP]', and '[PAD]' the numbers 101, 102, and 0, respectively. The other numbers are also assigned according to the pre-trained model’s values.

```

1 print(encoding['input_ids'][0])
2
3 tensor([ 101,   146,  1821,  4809,  1104,  1103,  1944,
4         1105,  1103,  4665,
5         117,   146,  1821,  4809,  1104,  1103,   182,
6         15626,  1105,  1103,
7         1514,   119,  2009,  1431,   146,  2215,   136,
8         1169,   170,  4105,
9         1200,  2640,  1518,  1435,  1106,  1143,  1303,
10        136,   152,   117,
11        1515,  1103, 10846,  1104,  4018,  1112,  1218,
12        1112,  1103, 10846,
13        1104,  2489,   117,  8640,  1122,  1136, 10228,
14        1191,   146,  6192,
15        1121,  1103,  1282,  1105,  1103,  7172,  1105,
16        1103,  2945,   136,
17        102,     0,     0,     0,     0,     0,     0,
18        0,     0,     0,
19         0,     0,     0,     0,     0,     0,     0,
20        0,     0,     0,
21         0,     0,     0,     0,     0,     0,     0,
22        0,     0,     0]])

```

**Listing 29:** Text Sample Numerical Representation

Finally, to close the encoding, we have the attention mask. In this mask the significant tokens receive a one and the padding tokens a zero. BERT will use this mask to give more relevance to the significant tokens in the sequence.

[illegible]

```
6 0, 0, 0, 0])
```

**Listing 30:** Attention Mask

## 6.2.4 BertClassifier

### 6.2.4.1 Bert Hyper-parameters

The BERT setup used for depression detection will now be briefly described.

```
1 RANDOM_SEED = 42
2 MAX_LEN = 50
3 BATCH_SIZE = 16
4 CLASSES = 2
```

**Listing 31:** Hyper-parameters for BERT

The seed number allows a model's stable execution over multiple training sessions. It initializes the weights of the neural networks randomly, but at the same initial values. The next parameter is the maximum length. This value will limit the length of text that the BERT model will process for each observation. Since the computing resources are limited, the text is truncated to max length. This truncation implies that the model has to detect the depression at the beginning of each text sample. The batch size will be the number of examples that the model will process simultaneously. Again, computing resources limits batch size range of values. Finally, the number of classes corresponds to 2, depressed with label 1, and not depressed with label 0.

We define the `BERTClassifier` for the implementation of the model. We derive this class from `nn.Module` (e.i Neuronal Network Module) from PyTorch. In initializing this class, we take the pre-trained BERT model, and add to it a Dropout and Linear layer. The Linear layer takes as arguments `self.bert.config.hidden_size` and `n_classes`, which correspond to the number of inputs and the number of outputs, respectively. The number of neurons at the output of BERT is equal to 768. This number will be the value of inputs for the Linear layer that we add to the pre-trained model for the classification. The `n_classes` corresponds to 2 because we are classifying two states, depressed or not. To reduce overfitting, i.e. the training and validation precision obtain a similar value, we have included a Dropout layer of 0.3. This layer will randomly turn off 30% of the neurons during the model's

training. The objective is a better generalization. In the second class' function, we define how these layers are connected. This function needs two inputs: the `inputs_ids` (the numerical representation of the tokens), and the `attention_mask`, which focuses BERT's attention. The first layer of our classifier will be the BERT model, which receives the `inputs_ids` and `attention_mask` as inputs. From the model outputs, we take `cls_output` that corresponds to the classification token. This output will go through the Dropout layer, which will finally be processed by the Linear layer. The Linear layer will give us the classification.

```

1 class BERTClassifier (nn.Module):
2
3     def __init__ (self, n_classes):
4         super(BERTClassifier, self).__init__()
5         self.bert = BertModel.from_pretrained(
6             PRE_TRAINED_MODEL_NAME)
7         self.drop = nn.Dropout(p=0.3)
8         self.linear = nn.Linear(self.bert.config.hidden_size ,
9             n_classes)
10
11     def forward( self, input_ids, attention_mask):
12         _, cls_output = self.bert(
13             input_ids = input_ids,
14             attention_mask = attention_mask)
15         drop_out = self.drop(cls_output)
16         output = self.linear(drop_out)
17         return output

```

**Listing 32:** Bert Classifier

#### 6.2.4.2 Bert Training and Validation Functions

For the training presented in this report, we used 5 epochs. We optimized our implementation with AdamW, an algorithm similar to the descending gradient. This optimizer receives the model parameters and the learning rate, also to avoid it modifying the model parameters, we select False for its `correct_bias` argument. The number of steps is equal to the number of batches delivered by the dataloader function multiplied by the number of epochs. The scheduler is responsible for reducing the learning range in each of the iterations. The scheduler arguments are the optimizer, the `num_warmup_steps` equal to zero because we want the learning rate to decrease from the first iteration, and the `num_training_steps`, which is equal to



the total of `total_steps`. Finally, we define the loss function, which we take from the `nn.CrossEntropyLoss()` module.

```

1 EPOCHS = 5
2 optimizer = AdamW(model.parameters(), lr=2e-5,
3                     correct_bias = False)
4 total_steps = len(train_data_loader)* EPOCHS
5 scheduler = get_linear_schedule_with_warmup(
6     optimizer,
7     num_warmup_steps= 0,
8     num_training_steps= total_steps)
9 loss_fn = nn.CrossEntropyLoss().to(device)

```

### Listing 33: Training

For training we define the `train_model` function. In this function we set the model on training mode, initialize the `losses` list to store the error. Also we assign zero to (`correct_predictions`) of zero. Then, for each batch that the dataloader delivers, the function sends the `input_ids`, the `attention_mask` labels to the GPU, which are then assigned as arguments to the model to obtain `outputs`. In `preds`, they correspond to the predictions calculated by BERT, the function uses `torch.max` to store the maximum value. The model has two outputs after being adjusted for depression detection. So, if the first output is maximum, it means that the person suffers from depression. If the second output is maximum, it means that the person does not have depression. After making the predictions, the function calculates the loss using the output and the original label. It also calculates the precision, which corresponds to the accumulated sum of the original labels' correct predictions. This data is stored in `losses`. Next, the function takes the value of the loss and executes backpropagation to update the weights. The function uses `nn.utils.clip_grad_norm(model.parameters(), max_norm = 1.0)` to cut the gradient value preventing the gradient from stagnating. Finally, the function updates the weights, the training rate, and reset the optimizer weights for the next batch process. The training function outputs the average value of the precision and the average value of the loss.

```

1 def train_model(model, data_loader, loss_fn, optimizer,
2                 device, scheduler, n_examples):
3     model = model.train()
4     losses = []
5     correct_predictions = 0
6     for batch in data_loader:

```

```

6         input_ids = batch['input_ids'].to(device)
7         attention_mask = batch['attention_mask'].to(device)
8         labels = batch['label'].to(device)
9         outputs = model(input_ids = input_ids, attention_mask
= attention_mask)
10        _, preds = torch.max(outputs, dim=1)
11        loss = loss_fn(outputs, labels)
12        correct_predictions += torch.sum(preds == labels)
13        losses.append(loss.item())
14        loss.backward()
15        nn.utils.clip_grad_norm(model.parameters(), max_norm
= 1.0)
16        optimizer.step()
17        scheduler.step()
18        optimizer.zero_grad()
19    return correct_predictions.double()/n_examples, np.mean(
losses)

```

**Listing 34:** Training Function

We define the `eval_model` function for the validation. This function has a similar structure to the `train_model` function. However, instead of initializing the model in training mode, the function takes the re-trained model and initializes it in evaluation mode. Furthermore, after processing each batch, the function does not modify the model weights; instead, it only returns the average value of the precision and the average value of the loss.

```

1 def eval_model(model, data_loader, loss_fn, device,
n_examples):
2     model = model.eval()
3     losses = []
4     correct_predictions = 0
5     with torch.no_grad():
6         for batch in data_loader:
7             input_ids = batch['input_ids'].to(device)
8             attention_mask = batch['attention_mask'].to(
device)
9             labels = batch['label'].to(device)
10            outputs = model(input_ids = input_ids,
attention_mask = attention_mask)
11            _, preds = torch.max(outputs, dim=1)
12            loss = loss_fn(outputs, labels)
13            correct_predictions += torch.sum(preds == labels)
14            losses.append(loss.item())

```

```

15     return correct_predictions.double()/n_examples, np.mean(
        losses)

```

**Listing 35:** Validation Function

### 6.2.4.3 Training and Results

As shown in the following snip of code, after five iteration with BERT we obtained an accuracy 94% during training and 71% on validation.

```

1 for epoch in range(EPOCHS):
2     print('Epoch {} from {}'.format(epoch+1, EPOCHS))
3     print('-----')
4     train_acc, train_loss = train_model(
5         model,
6         train_data_loader,
7         loss_fn,
8         optimizer,
9         device,
10        scheduler,
11        len(df_train) )
12    test_acc, test_loss = eval_model(
13        model,
14        test_data_loader,
15        loss_fn, device,
16        len(df_test))
17    print('Training: Loss {}, accuracy: {}'.format(train_loss
18        , train_acc))
19    print('Validation: Loss {}, accuracy: {}'.format(
20        test_loss, test_acc))

```

**Listing 36:** Training

```

1 Epoch 1 from 5
2 -----
3 Training: Loss 0.4908973435473551, accuracy:
4   0.7506917317708334
5 Validation: Loss 0.6028436693401696, accuracy:
6   0.7157389322916666
7 Epoch 2 from 5
8 -----
9 Training: Loss 0.36380973518741183, accuracy:
10  0.83880615234375
11 Validation: Loss 0.768554817569869, accuracy: 0.71435546875
12 Epoch 3 from 5

```

```

10 -----
11 Training: Loss 0.26638933907906903, accuracy:
    0.9001057942708334
12 Validation: Loss 1.0772863906270989, accuracy:
    0.7112630208333334
13 Epoch 4 from 5
14 -----
15 Training: Loss 0.21492615004043122, accuracy: 0.93115234375
16 Validation: Loss 1.1942933631016786, accuracy:
    0.7135416666666666
17 Epoch 5 from 5
18 -----
19 Training: Loss 0.19189651082607875, accuracy:
    0.9412638346354166
20 Validation: Loss 1.1942933631016786, accuracy:
    0.7135416666666666

```

**Listing 37:** Results

### 6.2.5 Insights gained

Collaborating in implementing the RNN and BERT models to detect depression has given us a satisfactory learning experience, especially for interdisciplinary work. During the implementation, we reinforced and expanded our knowledge of programming. We also learn and apply NLP concepts to understand how different written language representations are necessary for language processing. We recognized that those language representations have advantages for efficient processing, but they limit interpretation, which social sciences commonly require. Despite these limitations, we recognize interdisciplinary work's potential to create communication channels between experts from different disciplines. Thus we recommend more opportunities to develop interdisciplinary work.

## Listings

1	Parsing Folders . . . . .	7
2	Parsing XML files . . . . .	7
3	Tokenization . . . . .	9
4	Downsampling . . . . .	10
5	Vocab_to_integer . . . . .	11
6	Padding_features . . . . .	12
7	Tokenization example . . . . .	15
8	Encoding and Padding example . . . . .	15
9	Embedding Layer . . . . .	16
10	LSTM Layer . . . . .	16
11	hidden weights . . . . .	17
12	fully connected layer(Linear layer) . . . . .	18
13	Loss function . . . . .	19
14	Adam optimizer . . . . .	19
15	Model hyperparameters . . . . .	20
16	hidden state initialization . . . . .	20
17	Feedforward propagation . . . . .	21
18	Loss calculation and perform backpropagation . . . . .	21
19	Optimizer step . . . . .	21
20	Saving Model . . . . .	22
21	CV Hyperparameters . . . . .	23
22	Pretrained model . . . . .	25
23	Pre-trained Model . . . . .	26
24	Text Sample . . . . .	27
25	Tokens . . . . .	27
26	Tokens' numerical representation . . . . .	28
27	Encoding Functions . . . . .	29
28	Encoded Text Sample . . . . .	29
29	Text Sample Numerical Representation . . . . .	30
30	Attention Mask . . . . .	30
31	Hyper-parameters for BERT . . . . .	31
32	Bert Classifier . . . . .	32
33	Training . . . . .	33
34	Training Function . . . . .	33
35	Validation Function . . . . .	34
36	Training . . . . .	35

37	Results . . . . .	35
----	-------------------	----

## References

- [1] Jay Alammar. “The Illustrated BERT, ELMo, and co.(How NLP Cracked Transfer Learning)”. In: *Dec* 3 (2018), pp. 1–18.
- [2] Hayda Almeida, Antoine Briand, and Marie-Jean Meurs. “Detecting Early Risk of Depression from Social Media User-generated Content.” In: *CLEF (Working Notes)*. 2017.
- [3] *cross validation*, URL= [https://scikit-learn.org/stable/modules/cross\\_validation.html#cross-validation-evaluating-estimator-performance](https://scikit-learn.org/stable/modules/cross_validation.html#cross-validation-evaluating-estimator-performance), Accessed = 2020-12-16.
- [4] Jacob Devlin et al. “Bert: Pre-training of deep bidirectional transformers for language understanding”. In: *arXiv preprint arXiv:1810.04805* (2018).
- [5] Lina Gega, Isaac Marks, and David Mataix-Cols. “Computer-aided CBT self-help for anxiety and depressive disorders: Experience of a London clinic and future directions”. In: *Journal of clinical psychology* 60.2 (2004), pp. 147–157.
- [6] Sepp Hochreiter and Jürgen Schmidhuber. “Long short-term memory”. In: *Neural computation* 9.8 (1997), pp. 1735–1780.
- [7] Spencer L James et al. “Global, regional, and national incidence, prevalence, and years lived with disability for 354 diseases and injuries for 195 countries and territories, 1990–2017: a systematic analysis for the Global Burden of Disease Study 2017”. In: *The Lancet* 392.10159 (2018), pp. 1789–1858.
- [8] Simon Kemp. *Data-Reportal*. 2020. URL: <https://datareportal.com/reports/digital-2020-global-digital-overview1> (visited on 03/30/2020).
- [9] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).
- [10] Omer Levy and Yoav Goldberg. “Dependency-based word embeddings”. In: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 2014, pp. 302–308.
- [11] Isaac Marks and Kate Cavanagh. “Computer-aided psychological treatments: evolving issues”. In: *Annual review of clinical psychology* 5 (2009), pp. 121–141.

- [12] Health Quality Ontario. “Psychotherapy for Major Depressive Disorder and Generalized Anxiety Disorder: A Health Technology Assessment”. eng. In: *Ontario health technology assessment series* 17.15 (2017). PMC5709536[pmcid], pp. 1–167. ISSN: 1915-7398. URL: <https://pubmed.ncbi.nlm.nih.gov/29213344>.
- [13] *Overfitting isn't simple: Overfitting Re-explained with Priors, Biases, and No Free Lunch*, URL= <https://mlexplained.com/2018/04/24/overfitting-isnt-simple-overfitting-re-explained-with-priors-biases-and-no-free-lunch/>, Accessed = 2020-12-16.
- [14] Philip Resnik, Rebecca Resnik, and Margaret Mitchell. “Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality”. In: *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*. 2014.
- [15] Alex Sherstinsky. “Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm) network”. In: *Physica D: Nonlinear Phenomena* 404 (2020), p. 132306.
- [16] Nitish Srivastava et al. “Dropout: a simple way to prevent neural networks from overfitting”. In: *The journal of machine learning research* 15.1 (2014), pp. 1929–1958.
- [17] *Understanding LSTM Networks*, URL= <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>, Accessed = 2020-12-15.
- [18] Ashish Vaswani et al. “Attention is all you need”. In: *Advances in neural information processing systems* 30 (2017), pp. 5998–6008.
- [19] *What is dropout in neural networks?*, URL= <https://www.educative.io/edpresso/what-is-dropout-in-neural-networks>, Accessed = 2020-12-16.
- [20] WHO. *WHO, Depression*. 2020. URL: <https://www.who.int/news-room/fact-sheets/detail/depression>.