# Configuration Manual

MSc Research Project
Data Analytics

# Rohit Kumar Shrivas

Student ID: x19226403

School of Computing
National College of Ireland

Supervisor:     Hicham Rifai

**National College of Ireland**
**Project Submission Sheet**
**School of Computing**

| | |
|---|---|
| **Student Name:** | Rohit Kumar Shrivas |
| **Student ID:** | x19226403 |
| **Programme:** | Data Analytics |
| **Year:** | 2021 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Hicham Rifai |
| **Submission Due Date:** | 16/08/2021 |
| **Project Title:** | Configuration Manual |
| **Word Count:** | 621 |
| **Page Count:** | 8 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| | |
|---|---|
| **Signature:** | Rohit Kumar Shrivas |
| **Date:** | 16th August 2021 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

### Rohit Kumar Shrivas
### x19226403

## 1  Introduction

The main objective of this manual is to document the procedure and configuration of hardware utilized to provide a guideline to any user willing to produce the desired outcomes. This manual consists of code snippets used for exploratory data analysis, data preparation, model creation etc. alongside hardware and software specifications of the machine utilized for the implementation of this research project. The structure of this technical manual is as follows: Chapter 2, highlighting the hardware and software configuration requirements, Chapter 3, discusses the data collection technique adopted, chapter 4, sheds light on exploratory data analysis conducted over the data, and then chapter 5 discusses the implementation aspect of all the models.

## 2  Environment

### 2.1  Configuration of Hardware and Software

The hardware configuration of the machine utilized for the implementation of this research project is shown in Figure 1. It has Intel's i5 10th generation core with 2.11 GHz clock, 64-bit Windows 10 OS installed with 8 GB of onboard RAM.
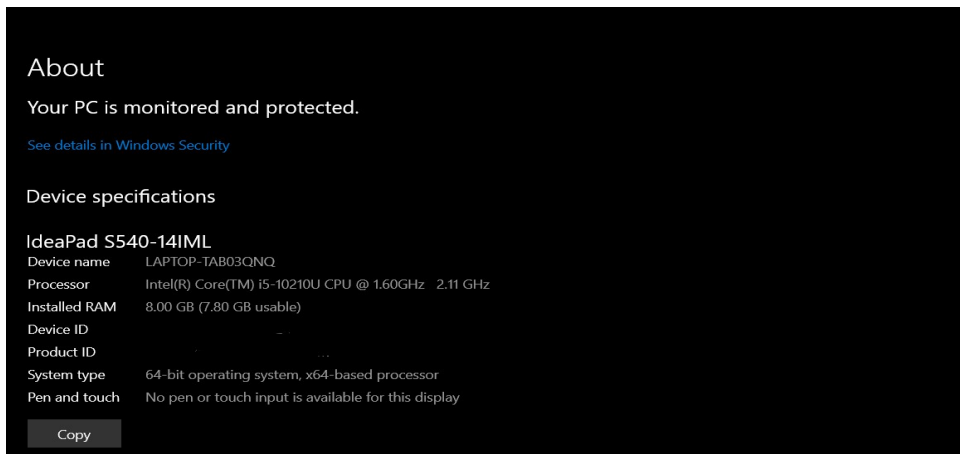


Figure 1: Hardware configuration of systemutilized for implementation of this research project

For the implementation of this project, Google Colaboratory'cloud based Jupyter notebooks have been utilized. For EDA purposes the default resource i.e., 12.69 GB of

RAM and 107.72 GB of Disk space is utilized. However, for RNN based Deep Learning Models, Graphics Processing Unit (GPU) accelerators have been utilized and for BERT based models, Tensor Processing Unit (TPU) have been utilized.

# 3  Collection of Data

The dataset has been downloaded from Kaggle which is an open platform. The link to the dataset is `https://www.kaggle.com/c/quora-insincere-questions-classification/data?select=train.csv`. This data is then uploaded into drive and the code is shown in Figure 2is utilized to upload it into Google's Colab instance.

```
1 # Authenticate and create the PyDrive client.
2 auth.authenticate_user()
3 gauth = GoogleAuth()
4 gauth.credentials = GoogleCredentials.get_application_default()
5 drive = GoogleDrive(gauth)

1 link_train = 'https://drive.google.com/file/d/1ktAT2mqPsRe2WmImqMIfFJkIvUwe_bMN/view?usp=sharing' # The shareable link

1 train_id = '1ktAT2mqPsRe2WmImqMIfFJkIvUwe_bMN'

1 downloaded = drive.CreateFile({'id':train_id})
2 downloaded.GetContentFile('train.csv')
3 train_df = pd.read_csv('train.csv')
```

Figure 2: Code for uploading data into Google Colaboratory instance from Drive

# 4  Exploratory Data Analysis

EDA was carried out to understand the flow of data. A pie chart has been plotted based on the classification category to understand the structure of data.

```
1 train_df['target'].value_counts().plot(kind = 'pie', labels = ['Sincere', 'Insincere'],
2     startangle = 90, autopct = '%1.0f%%')
```

Figure 3: Classification of data based on categories

Most frequent words found in both categories of data are then plotted along with wordcloud.

```
1  # convert huge strings to lists
2  Sincere = Sincere.split(' ')
3  Insincere = Insincere.split(' ')
4
5  # get the 15 most common words and their frequency values of each group using Counter
6  Sincere_count = Counter(Sincere).most_common(15)
7  Insincere_count = Counter(Insincere).most_common(15)
8
9  # function to generate input for plotting
10 def get_words_and_values(counter):
11     words = []
12     values = []
13     for i in range(0, len(counter)):
14         words.append(counter[i][0])
15         values.append(counter[i][1])
16     return words, values
17
18 # simple plot function
19 def freq_plot(words, values, tit = ''):
20     plt.bar(words, values)
21     plt.ylabel('Absolute Frequency')
22     plt.title(tit)
23     plt.xticks(rotation = 45)
24
25 # call the functions for SIncere questions first
26 words, values = get_words_and_values(Sincere_count)
27 freq_plot(words, values, tit = 'Frequency of the 15 most common words in Sincere questions')
```

Figure 4: Most frequent words appearing in sincere category



```
1  Sincere = train_df.loc[train_df['target'] == 0]
2  Insincere = train_df.loc[train_df['target'] == 1]
3
4  Sincere = ' '.join(question for question in Sincere['question_text'].astype(str))
5  Insincere = ' '.join(question for question in Insincere['question_text'].astype(str))
6
7  wordcloud = wc.WordCloud(max_font_size = 160, max_words = 150,
8                           background_color = 'black', width = 800,
9                           height = 500).generate(Sincere)
10 # plot it
11 plt.figure(figsize = (15, 10))
12 plt.imshow(wordcloud, interpolation = 'bilinear')
13 plt.axis("off")
14 plt.show()
```

Figure 5: Generating WordCloud of most frequent words appearing in sincere category

# 5 Implementation of Models

## 5.1 RNN based Deep Learning Models

Figure 6 shows the import of required libraries.

Figure 6: Libraries required by RNN based Deep Learning Models

Figure 7 shows the code utilized in splitting the data into train and validation set in the ratio of 80:20.



Figure 7: Splitting the data into train and validation

Figure 8 shows the code utilized for filling of '_NA_' values.



Figure 8: Filling '_NA_' values

Figure 9 shows the code utilized for tokenization of sentences.



Figure 9: Tokenizing the sentences

Figure 10 shows the code utilized for padding of sentences.



```
19 ## Pad the sentences
20 train_X = pad_sequences(train_X, maxlen=maxlen)
21 validation_X = pad_sequences(validation_X, maxlen=maxlen)
22
```

Figure 10: Padding the sentences

Figure 11 shows the code utilized for prediction of RNN based Deep learning Models.



```
1 predict_validation_y = model_4.predict([validation_X], batch_size=256, verbose=1)
2 for thresh in np.arange(0.1, 0.501, 0.01):
3     thresh = np.round(thresh, 2)
4     print("F1 score at threshold {0} is {1}".format(thresh, metrics.f1_score(validation_y, (predict_validation_y>thresh).astype(int))))
```

Figure 11: fnal classification prediction through RNN based Deep Learning models

## 5.2 Transformers based Models

Figure 12 shows the import of required libraries.



```
1 import re
2 import numpy as np
3 import pandas as pd
4 from collections import Counter
5 from afinn import Afinn
6 from nltk.corpus import sentiwordnet as swn
7 from nltk.corpus import stopwords
8 from nltk.stem import WordNetLemmatizer
9
10 from tokenizers import BertWordPieceTokenizer
11 from sklearn.model_selection import train_test_split
12 from sklearn import metrics
13
14 import tensorflow as tf
15 from tensorflow.keras.layers import Dense, Input
16 from tensorflow.keras.optimizers import Adam
17 from tensorflow.keras.models import Model
18 from tensorflow.keras.callbacks import ModelCheckpoint
19 import transformers
20 from tqdm.notebook import tqdm
21 from tokenizers import BertWordPieceTokenizer
22 import sentencepiece
```

Figure 12: Libraries required by Transformers based Models

Figure 13 shows the function created to fast tokenize the data.

```
[ ]    1 maxlen=192
       2 chunk_size=256
       3 def fast_encode(texts, tokenizer, chunk_size=chunk_size, max_length=maxlen):
       4     tokenizer.enable_truncation(max_length=maxlen)
       5     tokenizer.enable_padding(length=maxlen)
       6     all_ids = []
       7     #sliding window methodology
       8     for i in tqdm(range(0, len(texts), chunk_size)):
       9         text_chunk = texts[i:i+chunk_size].tolist()
      10         encs = tokenizer.encode_batch(text_chunk)
      11         all_ids.extend([enc.ids for enc in encs])
      12
      13     return np.array(all_ids)
```

Figure 13: Function for fast tokenization of data

Figure 14 shows the function created to build the transformer model.

```
 ▶   1 def build_model(transformer, max_len=512):
     2     input_word_ids = Input(shape=(max_len,), dtype=tf.int32, name="input_word_ids")
     3     #Replaced from the Embedding+LSTM/CoNN layers
     4     sequence_output = transformer(input_word_ids)[0]
     5     cls_token = sequence_output[:, 0, :]
     6     out = Dense(1, activation='sigmoid')(cls_token)
     7
     8     model = Model(inputs=input_word_ids, outputs=out)
     9     model.compile(Adam(learning_rate=1e-5), loss='binary_crossentropy', metrics=['accuracy'])
    10     model.summary()
    11     return model
```

Figure 14: Function for building transformer model

Figure 15 shows the code to detecting the no. of TPU clusters available.

```
1 try:
2     tpu = tf.distribute.cluster_resolver.TPUClusterResolver()
3     print('Running on TPU ', tpu.master())
4 except ValueError:
5     tpu = None
6
7 if tpu:
8     tf.config.experimental_connect_to_cluster(tpu)
9     tf.tpu.experimental.initialize_tpu_system(tpu)
10    strategy = tf.distribute.experimental.TPUStrategy(tpu)
11 else:
12    strategy = tf.distribute.get_strategy()
13
14 print("REPLICAS: ", strategy.num_replicas_in_sync)
```

Figure 15: Detecting th no. of TPU clusters available

Figure 16 shows the code which Tokenize the data using fast tokenizer function.

```
[ ]  1 #Tokenizing the samples
2
3 train_x = fast_encode(train_set['question_text'].astype(str), fast_tokenizer_distilbert, max_length=MAX_LEN)
4 val_x = fast_encode(validation_set['question_text'].astype(str), fast_tokenizer_distilbert, max_length=MAX_LEN)
5 train_y = train_set['target'].values
6 val_y = validation_set['target'].values
7 print(train_x.shape)
8 print(train_y.shape)
9 print(val_x.shape)
10 print(val_y.shape)
```

Figure 16: Function for fast tokenization of data

Figure 17 shows the code utilized in converting data into tensorflow compatible format.

```
[ ]    1 #Converting datasets in order to make it compataible with Tensorflow
       2
       3 train_dataset = (
       4     tf.data.Dataset
       5     .from_tensor_slices((train_x, train_y))
       6     .repeat()
       7     .shuffle(2048)
       8     .batch(BATCH_SIZE)
       9     .prefetch(AUTO)
      10 )
      11
      12 valid_dataset = (
      13     tf.data.Dataset
      14     .from_tensor_slices((val_x, val_y))
      15     .batch(BATCH_SIZE)
      16     .cache()
      17     .prefetch(AUTO)
      18 )
      19 print(train_dataset)
      20 print(valid_dataset)
```

Figure 17: Converting the data into tensorflow compatible format

Figure 18 shows the code utilized for prediction of Transformers based Models.

```
1 pred_val_y = DistilBERT_model.predict(val_x, batch_size=BATCH_SIZE)
2 for thresh in np.arange(0.1, .701, 0.01):
3     thresh = np.round(thresh, 2)
4     print("F1 score at threshold {0} is {1}".format(thresh, metrics.f1_score(val_y, (pred_val_y>thresh).astype(int))))
```

Figure 18: Final classification prediction through Transformers based models