

Configuration Manual

MSc Research Project Data Analytics

Shriya Gandhi Student ID: x19218079

School of Computing National College of Ireland

Supervisor: Jorge Basilio

National College of Ireland



MSc Project Submission Sheet

School of Computing

| Shriya Gandhi | |
|---------------|------|
| | |
| X19218079 | |
| | 2021 |

Student

Student ID:

Name:

| | Data Analytics | | 2021 |
|-------------------------|----------------------|---------|---------|
| Programme: | MSc Research Project | Year: | |
| Module: | Jorae Basilio | | |
| Lecturer: Submission | 16-08-2021 | | |
| Due Date: | | | |
| Project Title: | | nsier L | earning |

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.

<u>ALL</u> 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: | |
|------------|--|
| | |

Date:

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

Shriya Gandhi Student ID: x19218079

1 Introduction

This research project aims to automatically classify the insincere questions on Quora using the transformer-based models. In this configuration manual, we have outlined the steps to be followed for smooth replication of research project. The different steps of a project life cycle, starting from data collection until model testing and evaluation, are explained. Reference to code snippets is provided wherever required.

2 System Configuration

The entire project has been implemented on Google Colaboratory Pro with GPU (Tesla P100-PCIE-16GB) enabled. Google Colaboratory Pro version ensures that the model training session does not get disconnected because of RAM outage.

3 Dataset

The dataset for our research work is extracted from Kaggle.

4 Google Colaboratory Setup

The data extracted from Kaggle is uploaded to Google drive as in google colab we can access the data by simply mounting the drive. The code snippet for mounting of drive-in google colab can be found below.

Authorization for mounting google drive onto colab.



Google Drive successfully mounted.



Also, before proceeding towards coding, we first checked the GPU availability.

```
[1] # If there's a GPU available...
import torch
if torch.cuda.is_available():
    # Tell PyTorch to use the GPU.
    device = torch.device("cuda")
    print('There are %d GPU(s) available.' % torch.cuda.device_count())
    print('We will use the GPU:', torch.cuda.get_device_name(0))
    # If not...
else:
    print('No GPU available, using the CPU instead.')
    device = torch.device("cpu")
    There are 1 GPU(s) available.
    We will use the GPU: Tesla P100-PCIE-16GB
```

5 Installing Python Libraries for Data Pre-processing

The required libraries for data pre-processing are installed and imported.



```
[] train_df = pd.read_csv("/content/drive/MyDrive/Shriya/train.csv") #Read dataset
D
    train_df.head()
                           qid
                                                                question_text target
      0 00002165364db923c7e6
                                    How did Quebec nationalists see their province...
                                                                                      0
      1 000032939017120e6e44 Do you have an adopted dog, how would you enco...
                                                                                      0
        0000412ca6e4628ce2cf
                                    Why does velocity affect time? Does velocity a ...
                                                                                      0
      2
        000042bf85aa498cd78e How did Otto von Guericke used the Magdeburg h...
                                                                                      0
      3
        0000455dfa3e01eae3af
                                                                                      0
                                    Can I convert montra helicon D to a mountain b...
      Δ
[] train_df.shape
     (1306122, 3)
```

The uploaded dataset is read and stored inside a pandas dataframe.

6 Exploratory data analysis



```
Word cloud
```

```
#Wordcloud with commonly used words
comment_words = ''
   stopwords = set(STOPWORDS)
   # iterate through the csv file
   for val in train_df.question_text:
        # typecaste each val to string
       val = str(val)
       # split the value
       tokens = val.split()
       # Converts each token into lowercase
       for i in range(len(tokens)):
            tokens[i] = tokens[i].lower()
        comment_words += " ".join(tokens)+" "
   wordcloud = WordCloud(width = 800, height = 800,
                    background_color ='white',
                    stopwords = stopwords,
                    min_font_size = 10).generate(comment_words)
   # plot the WordCloud image
   plt.figure(figsize = (8, 8), facecolor = None)
   plt.imshow(wordcloud)
   plt.axis("off")
   plt.tight_layout(pad = 0)
   plt.show()
```

An additional column, quest_len, is added to python data frame. It stores the length of each question. The maximum, mean, median question lengths are identified.

```
[] #Max length of questions - sincere and insincere
train_df['quest_len'] = train_df['question_text'].apply(lambda x: len(str(x).split()))
print("Maximum length of a sincere question:", max(train_df[train_df['target']==0]['quest_len']))
maximum length of a sincere question: ", max(train_df[train_df['target']==1]['quest_len']))
Maximum length of a sincere question: 134
Maximum length of a insincere question: 64
[] #Mean length of questions
res = train_df['quest_len'].mean()
print("The mean length in words are : " + str(res))
The mean length of questions
res = train_df['quest_len'].median()
print("The median length in words are : " + str(res))
The median length in words are : " + str(res))
The median length in words are : " + str(res))
```

Using quest_len, plots to visualize the sincere and insincere question lengths are plotted.

```
#Plot for finding length of question tokens - sincere
sincere = train_df[train_df["target"] == 0]
plt.figure(figsize = (15, 8))
sns.distplot(sincere["quest_len"], hist = True, label = "sincere", color='blue')
plt.legend(fontsize = 10)
plt.title("Questions Length Distribution by Class", fontsize = 12)
plt.xlabel('Sincere Question Length')
plt.show()
```

```
#Plot for finding length of question tokens - insincere
insincere = train_df[train_df["target"] == 1]
plt.figure(figsize = (15, 8))
sns.distplot(insincere["quest_len"], hist = True, label = "insincere", color='red',norm_hist=True)
plt.legend(fontsize = 10)
plt.title("Questions Length Distribution by Class", fontsize = 12)
plt.xlabel('Insincere Question Length')
plt.show()
```

7 Data Pre-processing

The contractions such as didn't, couldn't present in the questions are expanded. Using lambda expressions, the steps of clean_text function are iterated over all the questions of the dataset. The pre-processed data is stored inside a new column, cleaned_questions.



Functions for removing punctuations and numbers are written and applied over the cleaned_questions text.



Finally, the data is validated for the presence of any duplicate records. The pre-processed file is then stored as a CSV file.



8 Model Implementation

We have implemented BERT and XLNET models for three samples of data. The steps for implementation of both models are similar. The below section briefs about the steps followed for XLNET with 10% data.

Installing the modeling dependencies



[] #importing libraries import pandas as pd import numpy as np from sklearn.model_selection import train_test_split

```
    #Importing required libraries for model training
from simpletransformers.classification import ClassificationModel
import pandas as pd
import logging
import sklearn
```

2021-08-14 21:12:31.452135: I tensorflow/stream_executor/platform/default/dso_loader.cc:53] Successfully opened dynamic library libcudart.so.11.0

The cleaned data is uploaded to google drive. The google drive is mounted in colaboratory, and data is read in pandas dataframe.

```
#loading drive
from google.colab import drive
     drive.mount('/content/drive')
    Mounted at /content/drive
⊡⇒
[ ] df=pd.read_csv("/content/drive/MyDrive/Shriya/clean_data.csv") #reading csv file
[ ] df=df.iloc[:,3:] #selecting the required column from dataframe
[] df.head()
         target
                                            cleaned questions
      0
              0
                   How did Quebec nationalists see their province...
      1
              0 Do you have an adopted dog how would you enco...
      2
             0
                     Why does velocity affect time Does velocity a ...
      3
              0 How did Otto von Guericke used the Magdeburg h...
      4
              0
                    Can I convert montra helicon D to a mountain b...
```

We have used the sample() function to extract 10% data from the original dataframe.

Note- For 50% and 70% data, the argument inside sample() changes to 0.5 and 0.7, respectively.

```
[ ] df=df.sample(frac=0.1) #sampling the data to select 10% records
[ ] df.shape
      (130533, 2)
```

The data is split into train-val-test at a ratio of 80-10-10 using the train_test_split function.

```
[] # splitting the data into training, test and eval dataset
    X = df['cleaned_questions']
    y = df['target']
[] #splitting in 80-10-10 ratio
    X_train, X_rem, y_train, y_rem = train_test_split(X,y, train_size=0.8)
    test_size = 0.5
    X_valid, X_test, y_valid, y_test = train_test_split(X_rem,y_rem, test_size=0.5)
[ ] print(X_train.shape), print(y_train.shape)
    print(X_valid.shape), print(y_valid.shape)
    print(X_test.shape), print(y_test.shape)
    (104426,)
    (104426,)
    (13053,)
    (13053,)
    (13054,)
    (13054,)
    (None, None)
```

Python's logging module is used to track the execution of model commands and errors if any.

```
[ ] logging.basicConfig(level=logging.INFO)
transformers_logger = logging.getLogger("transformers")
transformers_logger.setLevel(logging.WARNING)
```

After running multiple iterations, an optimum set of hyperparameters are chosen for our proposed models. We have used the Classification module provided by Simple Transformer distribution. The model, its tokenizer, maximum length of sequence, epoch, batch size,

learning rate, optimizer, dropout and more are given as arguments to the ClassificationModel function.

We have chosen evaluate_during_training to be True, thus in the train_model function, both training and evaluation data are passed.

The wandb_project stores the runtime data of the model. We have extracted the GPU utilization plot from wandb projects created for each run.

Note – For BERT models, the ClassificationModel arguments are bert and bert-base-uncased.



The eval_model is used to obtain the model performance results on evaluation data.



9 Predictions

Once the evaluation results look fine, predict function is used to test the model with test data. A confusion matrix and classification report are printed for each model.



The model is tested on few raw samples as well.

