

Configuration Manual

MSc Research Project Programme Name

Bhumika Sardana Student ID: x22105859@student.ncirl.ie

> School of Computing National College of Ireland

Supervisor: Vikas Tomer

National College of Ireland



MSc Project Submission Sheet

School of Computing

| Student Name | Bhumika Sardana |
|----------------------|---|
| Student ID | X22105959@student.ncirl.ie |
| Programme | Data Analytics |
| Year: | 2023 |
| Module: | Msc Research project |
| Supervisor: | Vikas Tomer |
| Submission Due Date: | 14/12/2023 |
| Project Title: | |
| | Sentimental analysis for Hinenglish-code mixed data |
| | using advanced word embeddings |
| Word Count: | 789 |
| Page Count: | 23 |

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 | Bhumika |
|-----------|---------|
| 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

Bhumika Sardana Student ID: x22105859@student.ncirl.ie

1 Introduction

This Configuration Manual lists together all prerequisites needed to duplicate the studies and its effects on a specific setting. A glimpse of the source for Exploratory Data analysis for is done followed by sentiment analysis, data preprocessing and cleaning and vectorization after that all the algorithms are created, and Evaluations is also supplied. After that recommendatory system is build and put together with the necessary hardware components as well as Software applications. The report is organized as follows, with details relating environment configuration provided in Section 2.

Information about data collection is detailed in Section 3. Exploratory Data Analysis is done in Section 4. Sentiment Analysis is included in Section 5. In section 6, the Data Clenaing is described. Section 7 provides details of Tokenisation. Details well about models that were created and tested are provided in Section 8. How the results are calculated and shown is described in Section 9.

2 System Requirements

The specific needs for hardware as well as software to put the research into use are detailed in this section.

2.1 Hardware Requirements

The necessary hardware specs are shown in Figure 1 below. MacOs M1 Chip, macOS 10.15.x (Catalilna) operating system, 8GB RAM, 256GB Storage, 24" Display.

| • • • | | | MacBook Pro | |
|---|--|--|-------------|--|
| Nardware | Hardware Oversiew | | | |
| Verlander Verlander | Hardware Overview: Model Name Chen Chen Manuelle | MacBook Pro Apple With Apple With Comparison of a difficiency 7466 A111 Difficulture VietworkLoops 960008103-000539840APT0018 Disabled | MacBook Pro | |
| Installations Language & Region Legacy Software Logs Managed Client Preference Panes Printer Software | | | | |
| Profiles Rew Support SmartCards Startur Jame | | | | |
| Startup ttems | | | | |

Figure 1: Hardware Requirements

2.2 Software Requirements

- Anaconda 3 (Version 4.8.0)
- Jupyter Notebook (Version 6.0.3)
- Python (Version 3.7.6)

2.3 Code Execution

The code can be run in jupyter notebook. The jupyter notebook comes with Anaconda 3, run the jupyter notebook from startup. This will open jupyter notebook in web browser. The web browser will show the folder structure of the system, move to the folder where the code file is located. Open the code file from the folder and to run the code, go to Kernel menu and Run all cells.

3 Data Gathering

The dataset is collected from https://github.com/l3cube-pune/code-mixednlp/tree/main. L3Cube-HingCorpus is the first large-scale real Hindi-English code mixed data in a Roman. It consists of 52.93M sentences and 1.04B tokens, scraped from Twitter. We also present HingBERT, HingMBERT, HingRoBERTa, and HingGPT. The BERT models have been pre-trained on codemixed HingCorpus.

4 Exploratory Data Analysis

Figure 2 includes a list of every Python library necessary to complete the project.



Figure 2: Necessary Python libraries

The Figure 3 represents the block of code to import data as text file nd enerate a list of sentence from the file.



The Figure 4 represents the block of code to convert list of sentences to data frame and printing data information and top 5 rows.



Figure 4: Data Information

In figure 5, the code to check total number of records and drop in case of missing data.

| print print print | :("Total no. of Data points:" ,len (df)) :("="*50) :(df[:10]) |
|--|---|
| Total | no. of Data points: 44453 |
| 0 sa 1 aa 2 gi 3 ek 4 ja 5 sa 7 ph 6 aa 9 aa | natan HI 0809 HI tiding EN luna EN in EN ham. p HI ne HI kaha HI tha HI apne HI video EN m. ll EN main HI first EN time EN bb13 HI ko HI. tHI nari HI ko HI duniya HI ke HI samne HI a. tti7 HI ki HI ho HI gea HI sadi HI stubborn . en EN i EN send EN this EN i EN mean EN this. nvee HI arre HI haan HI yaar HI agar HI ek H. ir HI tere HI katora HI khan EN ne HI kyu HI. ary HI ye HI adhikari HI apne HI power EN ka. HI gaya HI na HI apne HI idol EN wale HI fo. |
| df.dr df | ropna(inplace=True) |
| | Sentences |
| 0 | sanatan HI 0809 HI tiding EN luna EN in EN ham |
| 1 | aap HI ne HI kaha HI tha HI apne HI video EN m |
| 2 | gill EN main HI first EN time EN bb13 HI ko HI |
| 3 | ek HI nari HI ko HI duniya HI ke HI samne HI a |
| 4 | jatti7 HI ki HI ho HI gea HI sadi HI stubborn |
| | |
| 4444 | 8 arjuner HI going EN by EN thread EN jaakar HI |
| 4444 | 9 gill EN twitter EN pe HI mein HI apne HI free |
| 4445 | o sir HI aapne HI to HI sabko HI free EN ki HI a |
| 4445 | 1 999 EN rangnath EN aaj HI apne HI usually EN m |
| 4445 | 2 baat HI sahi HI hai HI but EN ab HI fozol HI k |
| 44453 | rows × 1 columns |

Figure 5: Missing Data

The Figure 6, checking train data.

| [] | trair | n = pd.read | _csv("/conte | ent/drive/MyDrive/Hingl: | ishTextSA/code-mixed-nlp-main/ | code-mixed-nlp-main/L3Cube- | -HingLID/train.txt", | sep='\t', names=['Wor | ds', 'Language']) |
|----|--|---|--|---|--|-----------------------------|----------------------|-----------------------|-------------------|
| [] | trair | n.info <mark>()</mark> | | | | | | | |
| | <clas Range Data # 0</clas | s 'pandas. Index: 968 columns (tr Column I Words 9 | core.frame.D 232 entries, otal 2 colum Non-Null Cou 968229 non-n | WataFrame'> 0 to 968231 mns): int Dtype uull object | | | | | |
| | 1 dtype memor | Language 9 s: object() y usage: 14 | 968232 non-n 2) 4.8+ MB | ull object | | | | | |
| [] | trair | n.head <mark>()</mark> | | | | | | | |
| | | Words Lan | guage | | | | | | |
| | | anatan | | | | | | | |
| | | 0809 | | | | | | | |
| | | tiding | | | | | | | |
| | | luna | EN | | | | | | |
| | 4 | | | | | | | | |
| 0 | print print print print | :("Total no :("Data poi :("Data poi :("="*50) :(train[:10 |). of Data po nts for Hind nts for Engl)) | pints:",len(train)) 1 words:", train.loc[tr Lish words:", train.loc | rain['Language']=='HI'].shape[[train['Language']=='EN'].shap | 0]) e[0]) | | | |
| Ð | Total Data Data | no. of Da points for points for | ta points: 9 Hindi words English wor | 968232 :: 693977 rds: 274255 | | | | | |
| | 0 sa 1 2 t 3 4 5 6 bc | Words Lang Inatan 0809 Liding luna in ham wuling | uage HI EN EN EN HI EN | | | | | | |
| | 7 8 | krte hai | HI HI | | | | | | |
| | 9 | аар | ні | | | | | | |

Figure 6: Train data

Figure 7 includes checking test data.

| [] | test = pd.read_csv("/content/drive/MyDrive/HinglishTextSA/code-mixed-nlp-main/code-mixed-nlp-main/L3Cube-HingLID/test.txt", sep='\t', names=['Words', 'Language']) |
|----|--|
| [] | test.info() |
| | <pre><class 'pandas.core.frame.dataframe'=""> RangeIndex: 193547 entries, 0 to 193546 Data columns (total 2 columns): # Column Non-Null Count Dtype 0 Nonds 193546 non-null object 1 Language 193547 non-null object dtypes: object(2) memory usage: 3.0+ MB</class></pre> |
| [] | test.head() |
| | WordsLanguage0hai1behan1behan2kiHi3IEN4hai |
| 0 | <pre>print("Total no. of Data points:",len(test)) print("Data points for Hindi words:", test.loc[test['Language']=='HI'].shape[0]) print("="\$90) print("="\$90) print(test[:10])</pre> |
| | Total no. of Data points: 193547 Data points for Hindi words: 136824 Data points for English words: 55723 |

Figure 7: Test data

The Figure 8 code segment to validation data.



Figure 8: Validation data

As seen in Figure 9, the list of all hindi words and generate wordcloud of them.



Figure 9: Hindi Wordcloud

As seen in Figure 10, the list of all english words and generate wordcloud of them.



Figure 10: English Wordcloud

5 Sentiment Analysis

The Figure 11, illustrate the sentiment analysis of sentence illustrate the code and to analyze value for positive or negative class based on final sentiment.



Figure 11: Sentiment Analysis

6 Data Cleaning

The Figure 12, illustrate the code to balance the data by applying undersampling.



Figure 12: Class Balancing

The Figure 13, illustrate the function to clear punctuations and contractions of the words in the sentence.



Figure 13: Clean punctuations and contractions

The Figure 14, illustrate the stopwords.



Figure 14: Stopwords

The Figure 15, illustrate basic variable initialized for text cleaning.

| <pre>#BASIC i=0 str1=' final_s all_pos all_neg s='' df</pre> | variables , tring=[] itive_words=[] # store words from +ve review ative_words=[] # store words from -ve review | vs here vs here. | | |
|--|--|---------------------|--|--|
| | | | | |
| | Sentences | Score | | |
| 1 | video video share link aksar poochti maulana s | 1 | | |
| 3 | nari duniya samne nagan sarir pradarshan sobha | 1 | | |
| 15 | sp party bade neta humara 1 | | | |
| 18 | malik bhut riyaat pmln gira dainge vote mangna 1 | | | |
| 19 | raffle debets karle dango karle chamcho karle 1 | | | |
| | | | | |
| 5533 | agr uthane pabandi ge biden kia bap leke aega | 0 | | |
| 7865 | shukla suar jakar suar idol post comment bheek | 0 | | |
| 29995 | ind tb uncle jagmohan kashmir governor questio | 0 | | |
| 13253 | block level transfer suru kijiye mantri ji shi | 0 | | |
| 5833 | meine sachi bhale twitter pr dekhen called | 0 | | |
| 31262 ro | 31262 rows × 2 columns | | | |

Figure 15: Variables

The Figure 16, illustrate the cleaning process checking each word in the sentence.



Figure 16: Cleaning Process

The Figure 17, illustrate the code to print clean sentence.

| df <mark>[</mark> "S€ | entences"] |
|-----------------------|--|
| 1 3 15 | video video share link aksar poochti maulana s nari duniya samne nagan sarir pradarshan sobha party bade neta humara |
| 18 | malik bhut riyaat pmln gira dainge vote mangna |
| 19 | raffle debets karle dango karle chamcho karle |
| 5533 | agr uthane pabandi biden kia bap leke aega tmhare |
| 7865 | shukla suar jakar suar idol post comment bheek |
| 29995 | ind uncle jagmohan kashmir governor question u |
| 13253 | block level transfer suru kijiye mantri shiksh |
| 5833 | meine sachi bhale twitter dekhen called |
| Name: | Sentences, Length: 31262, dtype: object |

Figure 17: Clean sentence

The Figure 18, illustrate wordcloud for positive words.



Figure 18: Wordcloud for Positive words

The Figure 19, illustrate wordcloud for negative words.



Figure 19: Wordcloud for Negative words

The Figure 20, the code to split data into training and test set.



Figure 20: Train test Split

7 Tokenisation

The Figure 21, illustrate the code for Bert Tokenizer.



Figure 21: Bert Tokenizer

The Figure 22, illustrate the code for TF-IDF Vectorizer



Figure 22: TfidfVectorizer

The Figure 23, illustrate the code for MURIL Tokenizer.





8 Machine Learning Models

8.1 LSTM Bert

```
lstm = Sequential()
lstm.add(Input(shape=(30,1)))
lstm.add(LSTM(64, activation='relu'))
lstm.add(Flatten())
lstm.add(Dense(32, activation='relu'))
lstm.add(BatchNormalization())
lstm.add(Dense(16, activation='relu'))
lstm.add(Dense(2, activation='relu'))
lstm.compile(optimizer='sgd', loss="categorical_crossentropy", metrics=["accuracy"])
lstm.summary()
Model: "sequential"
Layer (type)
            Output Shape
                                         Param #
lstm (LSTM)
                     (None, 64)
                                         16896
flatten (Flatten)
                     (None, 64)
dense (Dense)
                    (None, 32)
                                         2080
batch_normalization (BatchN (None, 32)
                                         128
ormalization)
dense_1 (Dense)
                    (None, 16)
                                          528
dense_2 (Dense)
                     (None, 2)
Total params: 19,666
Frainable params: 19,602
Non-trainable params: 64
lstm.fit(trainData, y_train )
<keras.callbacks.History at 0x7db6f386bc70>
```

Figure 24: Implementation of LSTM Bert

8.2 CNN Bert

| <pre>cnn = Sequential() cnn.add(Input(shape=(30,1))) cnn.add(ConvlD(128, 1)) cnn.add(MaxPooling1D(2)) cnn.add(Activation('relu')) cnn.add(ConvlD(64, 1)) cnn.add(MaxPooling1D(2)) cnn.add(Activation('relu')) cnn.add(ConvlD(32, 1)) cnn.add(MaxPooling1D(2)) cnn.add(Activation('tanh')) cnn.add(Flatten()) cnn.add(Dense(2, activation= cnn.compile(optimizer='adam') cnn.summary()</pre> | 'tanh')) , loss="categorical_crosse | ntropy", metrics=["accuracy"]) |
|---|--|-----------------------------------|
| Model: "sequential_1" | | |
| Layer (type) | Output Shape | Param # |
| conv1d (Conv1D) | (None, 30, 128) | 256 |
| max_pooling1d (MaxPooling1D) | (None, 15, 128) | ø |
| activation (Activation) | (None, 15, 128) | 0 |
| conv1d_1 (Conv1D) | (None, 15, 64) | 8256 |
| max_pooling1d_1 (MaxPooling 1D) | (None, 7, 64) | ø |
| activation_1 (Activation) | (None, 7, 64) | 0 |
| conv1d_2 (Conv1D) | (None, 7, 32) | 2080 |
| max_pooling1d_2 (MaxPooling 1D) | (None, 3, 32) | ø |
| activation_2 (Activation) | (None, 3, 32) | 0 |
| flatten_1 (Flatten) | (None, 96) | 0 |
| dense_3 (Dense) | (None, 2) | 194 |
| ====================================== | | |
| cnn.fit <mark>(</mark> trainData, y_train) | | |
| 929/929 [=================================== | | - loss: 0.7057 - accuracy: 0.5298 |

Figure 25: Implementation of CNN Bert

8.3 LSTM TFIDF

| <pre>lstm = Sequential() lstm.add(Input(shape=(n,1))) lstm.add(LSTM(12, activation='relu')) lstm.add(Platten()) lstm.add(Dense(32, activation='relu')) lstm.add(BatchNormalization()) lstm.add(Dense(2, activation='softmax')) lstm.compile(optimizer='adam', loss="categorical_crossentropy", metrics=["accuracy"]) lstm.summary()</pre> | | | | |
|---|---|--------------------------------------|--|--|
| Model: "sequential_2" | | | | |
| Layer (type) | Output Shape | Param # | | |
| lstm_1 (LSTM) | (None, 12) | 672 | | |
| flatten_2 (Flatten) | (None, 12) | | | |
| dense_4 (Dense) | (None, 32) | 416 | | |
| batch_normalization_1 (Batc hNormalization) | : (None, 32) | 128 | | |
| dense_5 (Dense) | (None, 2) | 66 | | |
| Total params: 1,282 Trainable params: 1,218 Non-trainable params: 64 | | | | |
| lstm.fit <mark>(</mark> X_train_Tfidf, y_tr | rain) | | | |
| 929/929 [=================================== | =======] - 10s 9ms/ste 0x7db6f27e2380> | ep - loss: 0.6873 - accuracy: 0.5306 | | |

Figure 26: Implementation of LSTM TF-IDF

8.4 CNN TFIDF

| <pre>cnn = Sequential() cnn.add(Input(shape=(n,1))) cnn.add(Conv1D(512, kernel_size=1, activation='relu')) cnn.add(Dense(8, activation='relu')) cnn.add(BatchNormalization()) cnn.add(Dense(2, activation='softmax')) cnn.compile(optimizer='Nadam', loss="categorical_crossentropy", metrics=["accuracy"]) cnn.summary()</pre> | | | | | |
|---|-----------------|-------------|--|--|--|
| Model: "sequential_3" | | | | | |
| Layer (type) | Output Shape | Param # | | | |
| conv1d_3 (Conv1D) | (None, 30, 512) | 1024 | | | |
| flatten_3 (Flatten) | (None, 15360) | 0 | | | |
| dense_6 (Dense) | (None, 8) | 122888 | | | |
| batch_normalization_2 (Batc hNormalization) | (None, 8) | 32 | | | |
| dense_7 (Dense) | (None, 2) | 18 | | | |
| | | ========= | | | |
| Total params: 123,962 Trainable params: 123,946 Non-trainable params: 16 | | | | | |
| | | | | | |
| <pre>cnn.fit(X_train_Tfidf, y_train)</pre> | | | | | |
| 929/929 [========================] - 5s 4ms/step - loss: 0.6619 - accuracy: 0.5482 <keras.callbacks.history 0x7db6f2d63c70="" at=""></keras.callbacks.history> | | | | | |

Figure 27: Implementation of CNN TFIDF

8.5 LSTM Muril

| <pre>lstm = Sequential() lstm.add(Input(shape=(n,1))) lstm.add(LSTM(64, activation='relu')) lstm.add(Flatten()) lstm.add(Dense(32, activation='relu')) lstm.add(BatchNormalization()) lstm.add(Dense(16, activation='relu')) lstm.add(Dense(2, activation='relu')) lstm.compile(optimizer='sgd', loss="categorical_crossentropy", metrics=["accuracy"]) lstm.summary()</pre> | | | | | | |
|--|---------------------|---------------------------|---------|------------------|----------------|------|
| Model: "sequential_4" | | | | | | |
| Layer (type) | Output | Shape | | Param # | | |
| lstm_2 (LSTM) | (None, | 64) | | ======= 16896 | | |
| flatten_4 (Flatten) | (None, | 64) | | | | |
| dense_8 (Dense) | (None, | 32) | | 2080 | | |
| batch_normalization_3 (Batc hNormalization) | (None | , 32) | | 128 | | |
| dense_9 (Dense) | (None, | 16) | | 528 | | |
| dense_10 (Dense) | (None, | 2) | | 34 | | |
| Total params: 19,666 Trainable params: 19,602 Non-trainable params: 64 | | | | | | |
| | | | | | | |
| lstm.fit(x_train_tokenized, | y_train | | | | | |
| 929/929 [===== <keras.callbacks.history (<="" at="" td=""><td>======= 0x7db6f:</td><td>====] - 6s 5r 29e6c20></td><td>ns/step</td><td>- loss: nan</td><td>- accuracy: 0.</td><td>4737</td></keras.callbacks.history> | ======= 0x7db6f: | ====] - 6s 5r 29e6c20> | ns/step | - loss: nan | - accuracy: 0. | 4737 |

Figure 28: Implementation of LSTM Muril

8.6 CNN Muril

| <pre>cnn = Sequential() cnn.add(Input(shape=(n,1))) cnn.add(Conv10(128, 1)) cnn.add(MaxPooling10(2)) cnn.add(Activation('relu'))) cnn.add(Activation('relu')) cnn.add(MaxPooling10(2)) cnn.add(Platten()) cnn.add(Platten()) cnn.add(Dense(2, activation= cnn.compile(optimizer='adam' cnn.summary()</pre> | 'tanh')) , loss= "catagorical_crosse | ntropy", metrics=["accuracy"]) | | |
|--|--|--------------------------------|--|--|
| Model: "sequential_5" | | | | |
| Layer (type) | Output Shape | Param # | | |
| conv1d_4 (Conv1D) | (None, 7, 128) | 256 | | |
| max_pooling1d_3 (MaxPooling 1D) | (None, 3, 128) | | | |
| activation_3 (Activation) | (None, 3, 128) | | | |
| conv1d_5 (Conv1D) | (None, 3, 64) | 8256 | | |
| max_pooling1d_4 (MaxPooling 1D) | (None, 1, 64) | | | |
| activation_4 (Activation) | (None, 1, 64) | | | |
| flatten_5 (Flatten) | (None, 64) | | | |
| dense_11 (Dense) | (None, 2) | 130 | | |
| Total params: 8,642 Trainable params: 8,642 Non-trainable params: 0 | | | | |
| cnn.fit(x_train_tokenized, y_train) | | | | |
| 929/929 [=================================== | | | | |

Figure 29: Implementation of CNN Muri

References

https://github.com/l3cube-pune/code-mixed-nlp/tree/main
Understanding TF-IDF (Term Frequency-Inverse Document Frequency) - GeeksforGeeks
MuRIL | A General Introduction to MuRIL - Analytics Vidhya
ganeshkharad/gk-hinglish-sentiment · Hugging Face
Understanding of LSTM Networks - GeeksforGeeks
Convolutional Neural Network (CNN) in Machine Learning - GeeksforGeeks