

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/342705915>

CATEGORIZATION AND TRANSLATION OPERATING SYSTEM'S ASSISTANCE IN EXPLICATION OF DIFFERENT BANGLADESHI ACCENTS

Article · June 2020

CITATIONS

0

READS

24

4 authors:



Nakib Aman Turzo
Varendra University

16 PUBLICATIONS 1 CITATION

[SEE PROFILE](#)



Pritom Sarker
Varendra University

10 PUBLICATIONS 1 CITATION

[SEE PROFILE](#)



Biplob Kumar
Varendra University

6 PUBLICATIONS 1 CITATION

[SEE PROFILE](#)



Niloy Kumar Shaha
University of Rajshahi

3 PUBLICATIONS 0 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



BD classic movie restoration Project [View project](#)



comparison [View project](#)

CATEGORIZATION AND TRANSLATION OPERATING SYSTEM'S ASSISTANCE IN EXPLICATION OF DIFFERENT BANGLADESHI ACCENTS

Nakib Aman Turzo

Lecturer,

Department of Computer Science & Engineering,

Varendra University,

Rajshahi, Bangladesh

Email: nakibaman@gmail.com

Pritom Sarker

B.Sc. in CSE,

Department of Computer Science & Engineering,

Varendra University,

Rajshahi, Bangladesh

Email: me.pritom@gmail.com

Biplob Kumar

B.Sc. in CSE,

Department of Computer Science & Engineering,

Varendra University,

Rajshahi, Bangladesh

Email: kumarbiplob336@gmail.com

Niloy Kumar Shaha

B.Sc. in CSE,

Department of Computer Science & Engineering,

Varendra University,

Rajshahi, Bangladesh

Email: niloyshaha20@gmail.com

ABSTRACT: National language of Bangladesh is Bengali and it's also the official language used frequently. Our paper's focal point was to categorize and differentiate West Bangla language or Bangladeshi Bangla accent in a Bengali sentence. We first amassed text from literature files. Then converted text sentence data to numeric data by using TF-IDF. After PCA application by MATLAB, final data set was being obtained. Our strategy for future will assist in developing an automatic software that detects if a sentence has been written in West Bangla or Bangladeshi Bangla and then it will do translation from one to another form. Differences between both Bangladeshi accents is already so minimum that only native speaker can identify them distinctively. There was no data available previously for this study. This work denoted that as if languages seems to be same but are unique and different in their own way and depicts the identity of two geographically separated regions. The major output of this work paid heed on identification of the form of language frequently used today. Many other studies could be conducted, based on the results of our study, on the effects of Sanskrit and Foreign literature

KEYWORDS: Bangladeshi Bangla, Inverse Data Frequency, Linear SVM, Principal Component Analysis, Python, Term Frequency, West Bangla

INTRODUCTION

Bangladesh's official and national language is Bengali with respect to Constitution's third article. 98% of Bangladeshis are fluent in Bengali as their first language. Bengali dialects are being classified in two dimensions i.e spoken vs. literary variations and prestige vs. regional variations. Spoken Bengali exhibits more variations than written one. Formal language including in speeches, news, announcements is in Cholit Bhasha. During Bengali standardization in late 19th and early 20th century, cultural elite mostly belong to regions like Kolkata, Hooghly, Howrah and Nadia. In both Bangladesh and west Bengal the standard today is based on West Central Dialect while the language has been standardized through to centuries of media and education with mostly speakers fluent in both their socio-geographical variety as well as the standard dialect used in the media. Differences in dialects are in three forms literary language vs. colloquial language, regional dialect vs. Standardized dialect and lexical variations. Dialectal names are originated from the districts where they belong. Standard form doesn't show much varieties across Bengali speaking areas of South Asia. Variations which are regional in spoken Bengali constitutes a dialect continuum.

Mostly speech differences occurs at a distance of few miles and have distinct forms among religious communitive vocabularies. Bengali Hindus tend to speak in Sanskritised Bengali while Bengali Muslims use Perso-Arabic. Western border dialects are spoken in the area known as Manbhum. There are many more minor dialects as well including those spoken in the bordering districts of Purnea and Singhbhum and among the tribals of eastern Bangladesh like Chakma and Hajong. Bengali's rich literature prior to 19th century was in rhymed verse. Writing system of Modern Bengali developed from an ancient Indian syllabary called Brahmi. Like all Brahmi scripts Bengali is being written from left to right with characters hanging from horizontal line. No distinction is present in upper and lower case letters.

LITERATURE REVIEW

Classifiers uncover contrasts in language yet not in cognizance. Cantonese use more than five sortal classifiers than Mandarin. 40% of things show up without classifier and 18% of Cantonese and 3% of Mandarin take a sortal [1].

Creation of a NP in Mandarin and Cantonese might be comprises of only a categorizer by using semantic measures to supersede their syntactic merchant [2].

Machine interpretation is a critical piece of Natural Language Processing (NLP) for transformation of one language to another. Interpretation comprises of language model, interpretation model and a decoder. A measurable machine interpretation framework was created to make an interpretation of English to Hindi. The model is created by utilizing programming in Linux condition [3].

Discourse and language preparing frameworks can be sorted by predefined etymological data use and is information driven and it utilized AI techniques to consequently concentrate and procedure applicable units of data are ordered as proper. In this way, a thought was misused utilizing ALISP (Automatic Language Independent Speech Processing) approach, with especially centering discourse handling [4]. Issue with numerous discourse understanding frameworks was the setting free language structure and enlarged expression structure syntaxes are requesting computationally.

Limited state language structures are effective however can't speak to the connection of sentence meaning. It was portrayed how language investigation can be firmly coupled by building up an APSG for examination of part and determining naturally. Utilizing this strategy proficient interpretation framework was manufactured that is quick contrasted with others [5]. In another exploration the mix of regular language and discourse preparing in Phi DM-Dialog and its cost-based plan of equivocalness goals were talked about. The synchronous understanding ability was made conceivable by a steady parsing and age calculation [6].

Change of language is the hardest assignment and a contextual investigation was accomplished for this exchange off. This remembered interpretation of customer's framework for restrictive language into programming dialects. Various components were considered that influence robotization level of language transformation [7].

In 1996 CJK Dictionary Publishing Society began an analytical task for the issues top to bottom and for making an elaborative streamlined Chinese and conventional Chinese information base with 100% exactness by working together with Basis Technology in creating advanced division [8].

In hardly any investigations discourse to content of words transformation were accomplished for reconciliation of individuals with hearing weaknesses. Improvement of programming was to help individual through rightness of elocution utilizing English phonetics. This product helps in acknowledgment of potential in English hearing [9].

A presentation of nonexclusive technique for changing over a composed Egyptian everyday sentence to diacritized Modern Standard Arabic (MSA) sentence which could without much of a stretch be reached out to be applied to different vernaculars of Arabic which could undoubtedly be applied to different lingos. A lexical obtaining of informal Arabic was done which is utilized to change over composed Egyptian Arabic to MSA [10]. A framework was likewise evolved in such manner which perceives two speakers in every one of Spanish and English and was constrained o 400 words. Discourse acknowledgment and language examination are firmly coupled by utilizing a similar language model [11].

In an examination by utilizing neural system transformation of content written in Hindi to discourse was done which has numerous applications in everyday life for daze. It is likewise utilized for teaching understudies. The report containing Hindi was utilized as information and neural system was utilized for character acknowledgment [12].

There is limitation of syntactic blunders in inconstancy and capacity in verifiable times of English. In nineteenth and twentieth century they become increasingly beneficial joined by significant expansions in capacity, variations and scope of lexical affiliation [13].

For transformation of Hindi content to discourse in Java Swings a Graphical User Interface has been planned on the grounds that it comprises of various dialects spoken in various zones [14].

As of late advances were made in discourse union has delivered synthesizers with exceptionally high comprehensibility yet the expectation and sound quality is as yet an issue. Be that as it may, its quality has arrived at a sufficient level for some applications [15].

There are numerous looks into likewise focused on acknowledgment exactness of discourse with installed spelled letter groupings. Various techniques got proposed to limit spelled letter portions and rename them with a particular letter recognizer [16].

Improvement report was set up for interpreter programming which mostly counterbalances the nonattendance of instructive devices that conference weakened, requirement for correspondence. For creating composed language abilities this apparatus could be utilized [17].

For changing over words into triplets Software framework changes over among graphemes and phonemes utilizing vocabulary based, rule based and information driven strategies. A shotgun incorporate these strategies in a half and half framework and includes etymological and instructive data about phonemes and graphemes [18].

An online discourse to content motor was produced for move of discourse into composed language continuously and it required exceptional procedures [19].

Examination of interpretation situations was done in subjective research. Vehicle of composed and communicated in language was fundamentally tested by considering the ramifications of comparable issues. Interpretation as essential issue and how its managed issues raised by portrayal that would be worry for all analysts [20].

A similar work was also performed on differentiation and translation of Sadhu and Cholit language which was the basis of inter-conversion of other languages. Here Linear Discriminant Analysis performed best and speed prediction was also done. That is why Sadhu didn't remain complex language [21].

METHODOLOGY

Sum total of 28550 sentences were taken into account for this task. Altogether 10800 Bengali sentences from 8 distinctive literatures were being selected and 17750 Bengali sentences from West Bengal were picked from 10 distinct literatures.

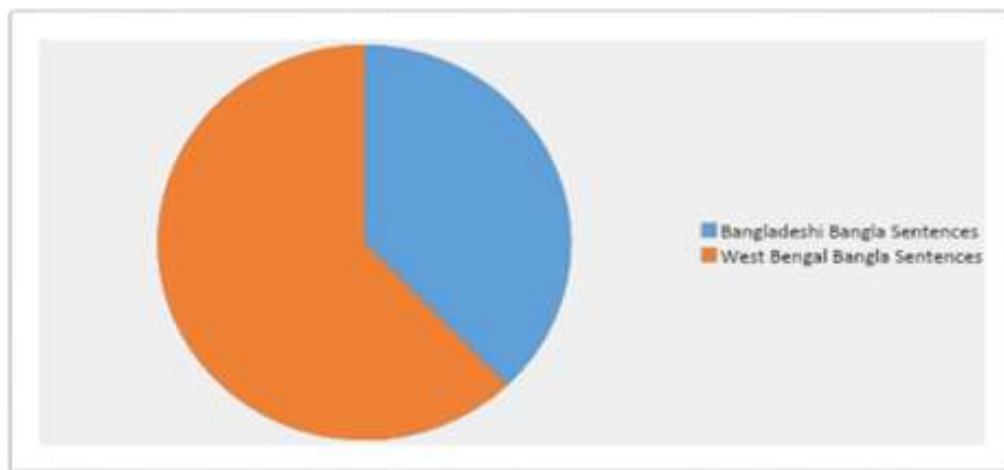


Fig 1

The methodological steps used are as follows:

First we amassed a .txt file literature and then got well defined sentences from the literature. From each of the sentence we conjectured stop word. Then text sentence data

is being metamorphosed to numeric data by utilizing TF-IDF. Final data set is obtained by application of PCA on data by using MATLAB and Python a variety of machine learning algorithms on the information set. At the ending point through analytical approach inspection is being done.

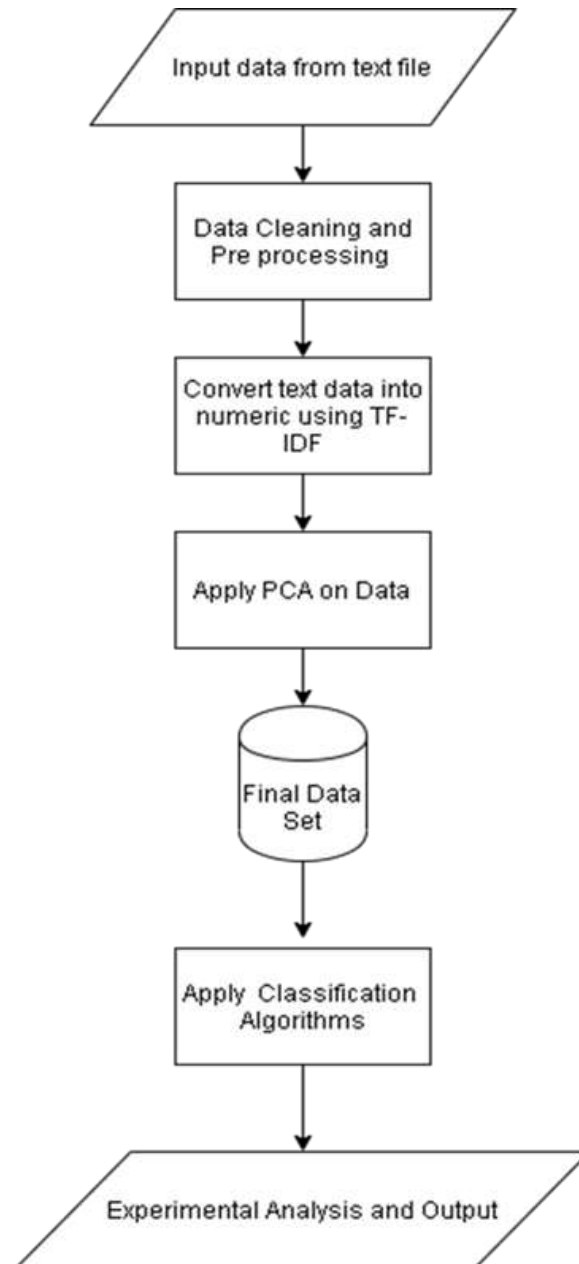


Fig 2: Work Flow

Data Clean

We have non-English (which got filtered out before or after processing of natural language data) in our set of information. All the non-English words got axed from it by us. Natural Language Toolkit (NLTK) information center of python is being used for this purpose. We have all of the sentences in non-English in our information set. Ergo, after the moping through the process, on the norm we got 1983 data set. As far as

numeric categorization is concerned Sadhu is dubbed as numeric 0 and cholit is categorized as numeric 1.

Term Frequency–Inverse Document Frequency

An analytical statistic is a numerical or scientific form of statistic which is being contemplated to mirror the principal of word in a docket or corpus and is called Short Term Frequency-Inverse Document Frequency (TF-IDF). This factor has weightage in retrieving information, text mining and user modeling through hunting of this data.

Term Frequency (TF)

Frequency of a word which pops up in a docket divided by the gross number of words in the document. Every document has its own term frequency.

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{i,j}}$$

Inverse Data Frequency (IDF)

The log of the documents number divided by word w containing documents. Inverse data frequency determines the weight of rare words across all documents in the corpus.

$$idf(w) = \log\left(\frac{N}{df_t}\right)$$

TF-IDF is simply the TF multiplied by IDF

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

Our most work is being done from Scikit-Learn which is TF-IDF Vectorizer's class. Our text data is taken by it and converted to numeric information set. After this conversion, our data has 3394 features. We have so many less important features we can do features extraction using PCA.

Principal Component Analysis

For pages other than the first page, start at the top of the page, and continue in double-column format. The two columns on the last page should be as close to equal length as possible. A new coordinate system is being metamorphosed from data through orthogonal linear transformation so that each coordinate has greatest variance by scalar projection of data in an ordered way and so on. This is called principal component analysis. Principal component analysis is a class of Scikit-learn. Higher variance comes to lie in first coordinate which is called first principal component and the lower variance in second coordinate. Our information set has 1678 traits after application of principal component analysis. When applications of dimensions of principal component analysis got reduced and the data quality got lost.

In case of principal quality analysis, 95% caliber of data was being maintained. 95% of the quality of real data was preserved by setting value of 'n' components as 0.95. Our latest data has 1678 characteristics after application of principal component analysis.

10757	-0.01192	4.34E-05	8.09E-05	-0.00123	-0.00256	-0.0016	0.001695	-0.00034	0.007713	0.00146	0.010956	0.007934	0.004034	0.004848	-0.01433	-0.00748	-0.00198	-0.00073	0.001574	-0.00671	1
10758	-0.000208	0.000217	0.000625	0.000272	0.000585	-8.80E-05	0.001269	0.0009	0.001517	0.002014	-0.00112	-0.00039	0.000929	0.000603	0.000493	0.000257	0.001002	4.30E-05	8.37E-05	-0.00152	1
10759	9.21E-05	7.69E-05	8.23E-05	0.000904	-0.00012	-0.00028	8.79E-05	-0.00026	-0.00077	0.000141	-0.00041	0.000597	0.000384	0.000552	0.000143	-0.00113	-0.00094	-0.0005	0.000231	0.000619	1
10760	-0.000718	-0.0002	0.000345	-0.00054	-0.00127	-0.00179	0.000675	-0.00012	0.004179	0.000884	0.002102	0.005528	0.002564	0.002548	-0.00043	-0.00055	-6.94E-05	0.000301	0.000149	-0.00288	1
10761	0.001698	-0.00099	-0.00172	-0.00393	-0.00011	-0.00048	0.002739	-0.00336	-0.00179	-0.00161	-0.02471	0.001913	0.010012	0.004037	0.01329	0.000315	-0.04341	0.074356	0.020738	-0.01918	1
10762	-0.00085	0.000466	0.000174	-0.00055	-0.00028	-0.00072	0.000628	-0.00078	-0.00116	-0.00024	0.00052	0.000409	0.000366	-0.00052	0.001375	-0.00034	-0.00041	0.000528	0.001636	0.000102	1
10763	0.001362	-0.00087	-0.00178	-0.00305	1.86E-05	-0.00734	0.003492	-0.00321	-0.00116	-0.00146	-0.02377	0.003076	0.009615	0.004827	0.012243	0.00454	-0.04334	0.072917	0.020645	-0.01945	1
10764	0.00078	0.000148	-0.00025	-0.0001	-0.00017	0.000524	-0.0002	0.000226	1.88E-05	-1.67E-05	-0.00025	0.000359	0.000536	0.00013	-0.00106	0.000658	-0.00069	0.000793	0.001469	-0.0005	1
10765	-0.00115	0.000243	-2.91E-05	0.00031	0.001741	-0.00119	-0.00011	-4.66E-05	3.86E-05	-0.00038	0.000179	0.000628	-3.37E-05	9.70E-05	0.000789	0.001083	-0.00082	-0.00038	0.000387	0.001045	1
10766	-0.00023	-0.00069	0.001039	-0.00065	0.003152	-0.00095	0.002139	-2.08E-05	-0.00372	-0.00073	0.001617	-0.00598	0.009424	-0.0028	-0.00199	0.004351	0.001574	0.002968	-0.00104	0.001998	1
10767	-0.00012	-0.00096	-0.00082	-0.00022	-0.00012	0.000634	-0.00042	-0.00025	0.000222	0.000151	-0.0002	-0.00024	0.000726	-0.00029	-0.00014	0.000264	-0.00027	0.000117	-0.00106	0.000155	1
10768	-0.00012	0.000189	7.42E-05	0.000922	0.001217	0.000924	0.00409	0.00188	0.002778	0.004995	-0.00385	-0.00287	-0.00044	-4.40E-05	-0.00087	0.00024	-0.00065	-0.00052	0.000721	0.003862	1
10769	-0.00005	0.000693	-0.00023	-0.00039	-0.00034	0.001592	0.0018	-0.00151	-0.00282	-0.00129	-0.00022	0.000265	-0.00094	0.0012	0.001381	0.001258	-0.00027	0.001082	-0.00119	0.00103	1
10770	0.000951	0.000511	9.31E-05	0.000449	-5.92E-05	0.001233	-0.00128	0.000909	0.00045	0.000927	-0.00187	0.001567	-0.00045	-0.00026	0.001827	-0.00014	-0.00273	0.002569	0.001713	-0.00012	1
10800	0.000229	-3.08E-05	-0.00029	0.000175	-0.00042	-0.00105	0.000815	-0.00038	-0.00028	-0.00028	0.000927	-0.00071	0.000371	-0.00054	0.000814	0.00165	-0.00168	-0.00017	0.001289	1	
10801	0.00154	-0.00016	0.00104	-0.00101	0.001012	-2.49E-05	0.000378	0.000416	-0.00017	-0.00114	0.000535	-0.00028	0.00085	0.00017	0.000153	0.000121	0.000933	0.001054	-0.00027	0.000439	1
10802	-0.0006	-7.53E-05	-0.00044	-0.00013	-0.00018	-9.30E-05	3.34E-06	-0.00051	0.000438	0.000295	-0.00031	-0.00051	8.37E-05	7.83E-05	-0.00059	-6.51E-05	-3.08E-05	-0.0004	-0.00029	-9.91E-05	1
10803	0.00314	-0.00111	0.000193	-0.00075	0.001266	0.000372	-1.18E-05	-3.15E-05	-0.00019	0.000431	-0.00058	0.000454	-1.01E-05	-9.86E-05	8.96E-05	-0.00034	0.000241	-0.00022	0.000369	-0.00024	0
10804	-0.00012	9.03E-05	-0.00022	-0.0002	-3.64E-05	-0.00034	8.07E-05	-0.00019	0.000431	-0.00058	0.000454	-1.01E-05	-9.86E-05	8.96E-05	-0.00034	0.000241	-0.00022	0.000369	-0.00024	-0.00029	0
10805	-0.00043	8.41E-05	-0.0008	0.000268	-0.00041	-0.00278	-0.0001	-0.0002	0.000705	0.00101	0.000252	0.000867	0.000832	0.000349	-0.00071	-0.00023	-0.00051	-0.00075	0.00074	-0.00029	0
10806	0.000317	8.21E-05	-0.00031	-0.00035	0.00101	-0.00022	-0.00033	0.000216	9.46E-06	0.000585	0.000643	-0.00038	-0.00027	0.000702	-0.00085	-0.00082	-0.00078	0.00112	-0.00036	-0.00119	0
10807	0.001	-0.00061	0.000193	0.000763	-0.00139	2.33E-05	-0.00025	0.000245	0.000158	0.000488	-0.00037	0.000401	-0.00056	-0.00015	-0.00037	-0.00032	0.000258	-0.00028	-3.21E-05	-0.00127	0
10808	-4.90E-06	-0.00036	8.61E-05	-0.00042	-0.00079	-0.00141	-0.00023	0.00013	-0.00058	0.000285	0.001045	-0.00038	-0.00336	0.001255	-0.00293	0.000491	0.001405	2.66E-05	0.001313	0	
10809	-0.00035	0.000129	-0.00014	0.000374	-0.00127	-0.0005	0.000225	-0.0006	-0.0004	0.000241	5.31E-05	-0.0002	2.06E-05	-8.06E-05	-0.00043	-0.0003	-0.00065	0.000123	-8.27E-05	0	
10810	-0.00084	0.000853	-0.00226	-0.0021	-0.00197	-0.00265	-0.0013	-0.00064	0.000208	0.000428	0.001018	0.002158	0.000782	0.001298	0.000562	0.001458	-7.55E-05	0.000987	-0.00117	0.000422	0
10811	-0.0004	-0.0004	0.000357	-0.00144	-0.00034	0.000212	0.00112	-0.00074	-0.0012	-0.00044	0.001471	0.000204	0.001365	0.000907	-0.00137	-0.0008	-0.00107	-0.00028	-0.00055	0.000828	0
10812	-0.00053	0.000234	-0.00082	0.00054	0.00096	0.000342	-0.00027	-0.0005	-1.15E-05	8.12E-05	0.000976	7.50E-05	-0.00038	0.00083	-0.00134	-0.00064	-0.00018	-0.00043	-0.00037	-0.00133	0
10813	0.000146	-0.00012	-7.13E-05	0.000285	4.48E-05	-0.00092	-0.0004	-0.00016	5.76E-05	-0.00035	0.000358	6.82E-05	9.29E-05	-0.00022	-0.00027	0.000883	-0.00047	-0.00063	0.000163	-0.00051	0
10814	3.86E-06	-0.00082	0.000263	-0.00066	0.000988	-6.66E-05	0.000189	0.000124	0.000988	0.00091	0.002846	0.002228	0.001244	0.000288	-0.00026	-0.00266	0.000299	-2.66E-05	0.000408	-0.0008	0
10815	2.97E-05	0.000101	0.000483	-8.80E-05	-0.0004	-1.63E-05	0.000808	-0.00026	-0.00056	0.000117	-5.01E-05	-0.00036	0.000484	-7.05E-05	-9.13E-05	0.00054	0.000585	-0.00011	2.46E-05	-0.00075	0
10816	-8.28E-06	-0.00049	0.000332	-0.00019	0.000254	-0.00033	-0.00019	-2.41E-05	-2.77E-05	-0.00019	0.000764	0.000185	0.000822	0.000359	-0.00039	0.000867	0.000647	-3.78E-05	0		

Fig 3

There are 1194 different sides of numeric data in which the last field denotes 1 for Bangla language of Bangladesh while 0 for Bangla language of West Bengal.

RESULTS AND EXPERIMENTAL ANALYSIS

After implementing dataset in MATLAB results and factors for total misclassification of top 4 classifiers are as follows:

Table 1

Classifier Name	Prediction speed	Training Time	Total Misclassification Cost	Accuracy
Linear SVM	1200	2715.7	3385	76.3%
Quadratic SVM	46	3564.7	3324	76.7%
Medium Gaussian SVM	58	1775.7	3448	75.8%
Bagged Trees	1600	740.01	4041	71.7%
Subspace Discriminant	110	1057.9	3364	76.4%

Bagged Trees depiction through prediction speed graph has maximum speed of anticipation while the Quadratic SVM has the lowest anticipation speed. Categorizers like Naïve Bayes and Decision Tree are utilized but ameliorated due to less precision.

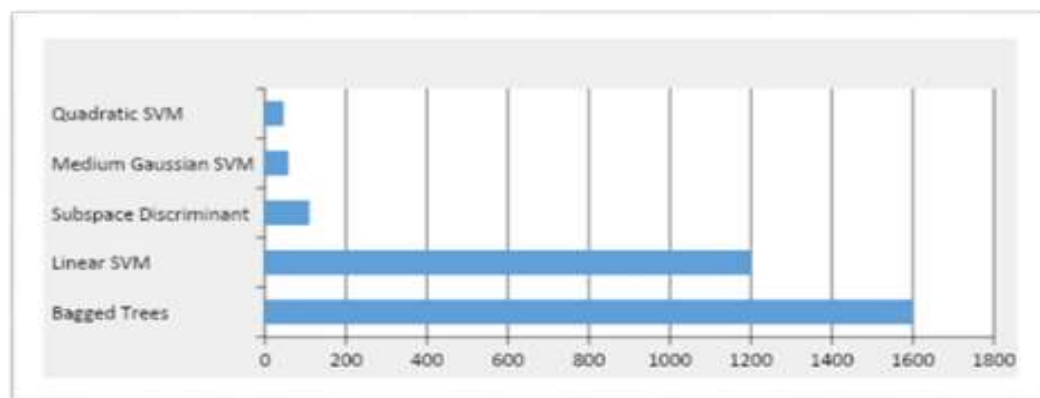


Fig 4

Subspace Discriminant comes after Bagged Tree which has highest training time. SVM categorizers are much slow in this regard.

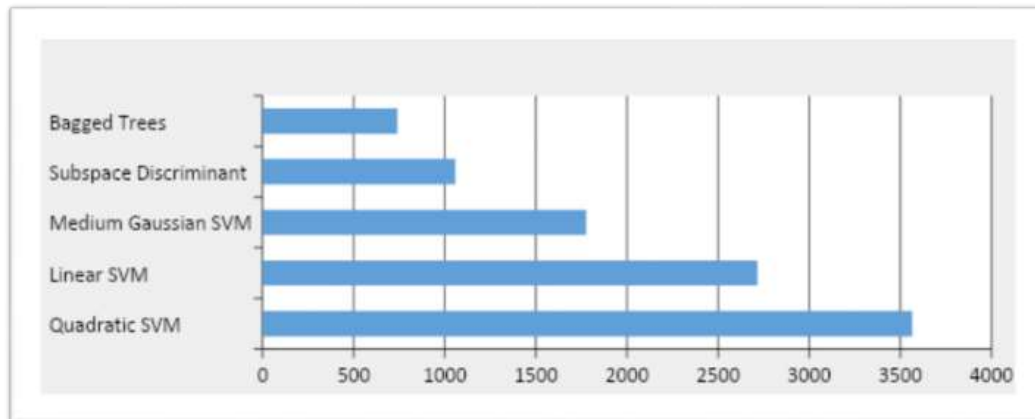


Fig 5

Fascinating features depicted here is total misclassification cost. Bagged tree which has previously highest value is the only categorizer lowest here while the others are similar.

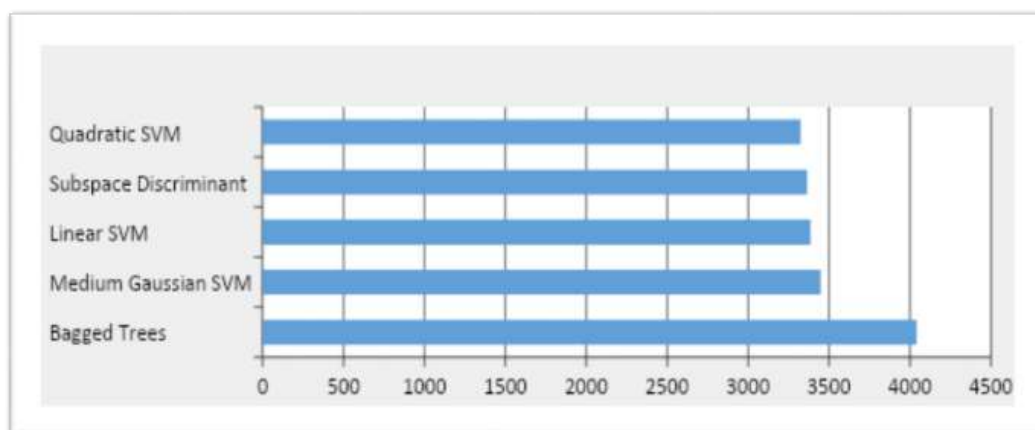


Fig 6

Quadratic SVM gave the highest accuracy followed by Subspace Discriminant.

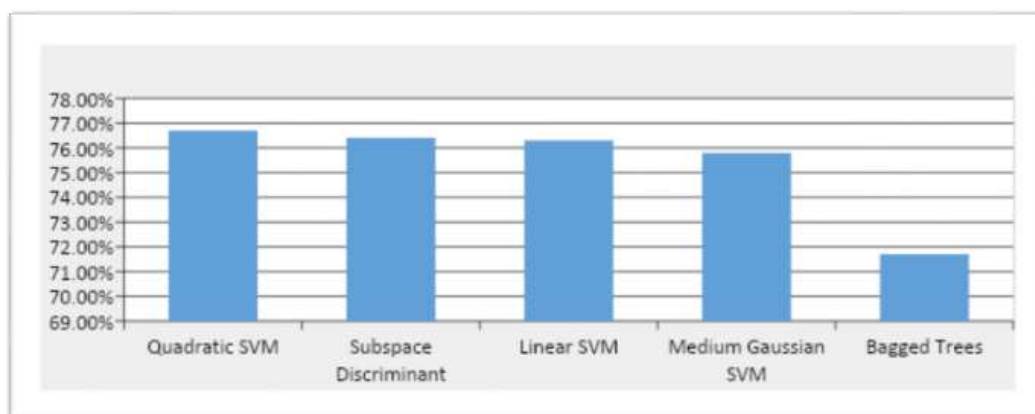


Fig 7

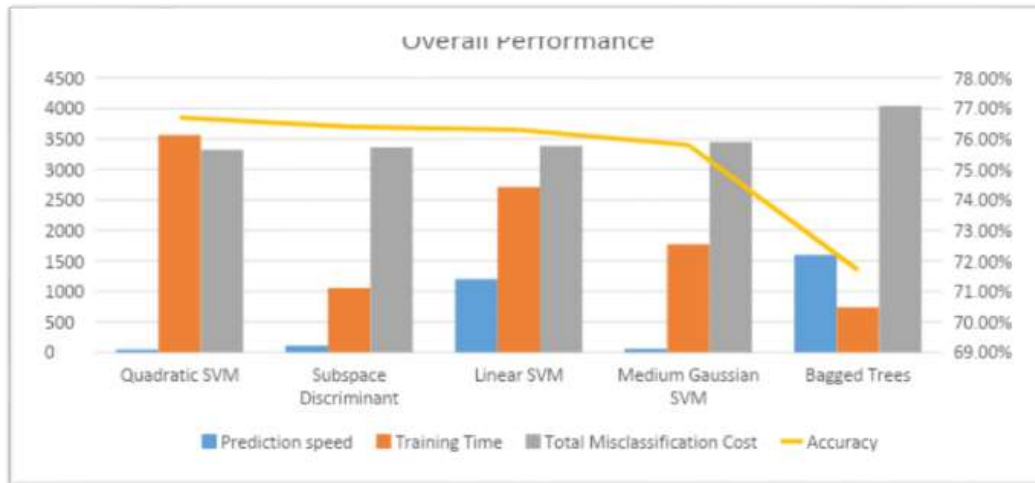


Fig 8

It can be observed from the functions that Quadratic SVM is highly precise and has high training time. Linear SVM is at 3rd rank with high precision. So after optimization in MATLAB, Linear SVM is beneficial for categorizing Bangladeshi Bangla and West Bengali Bangla.

Here are the preset values used for the classifiers –

Table 2

Quadratic SVM	Linear SVM	Medium Gaussian SVM	Boosted Trees	Subspace Discriminant
Model Type	Model Type	Model Type	Model Type	Model Type
Preset: Quadratic SVM	Preset: Linear SVM	Preset: Gaussian SVM	Preset: Bagged Trees	Preset: Subspace Discriminate Kernel
Kernel function: Quadratic	Kernel function: Linear	Kernel function: Gaussian	Kernel function: Bag	Kernel function: Subspace
Kernel scale: Automatic	Kernel scale: Automatic	Kernel scale: 35	Learner Type: Decision tree	Learner Type: Discriminant
Box constraint level: 1	Box constraint level: 1	Box constraint level: 1	Maximum number of split: 28551	Maximum number of split: 30
Multiclass method: One-vs-One	Multiclass method: One-vs-One	Multiclass method: One-vs-One	Number of learner: 30	Number of learner: 597
Standardize data: true	Standardize data: true	Standardize data: true		

Confusion Matrix of 5 classifiers:

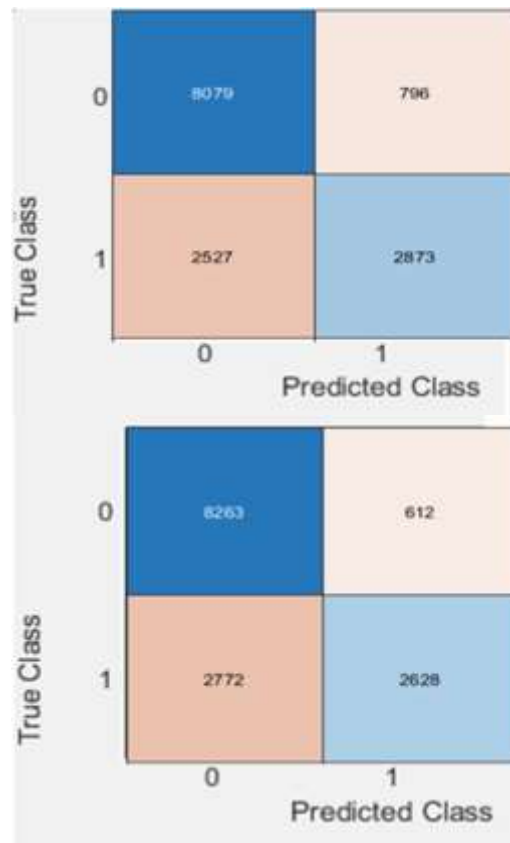


Fig 9: Quadratic SVM

Fig 10: Linear SVM

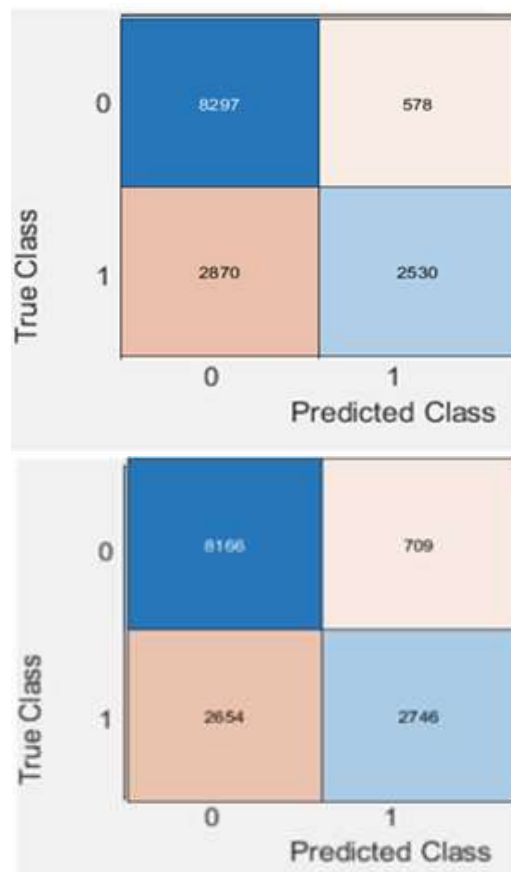


Fig 11: Medium Gaussian SVM

Fig 12: Boosted Trees

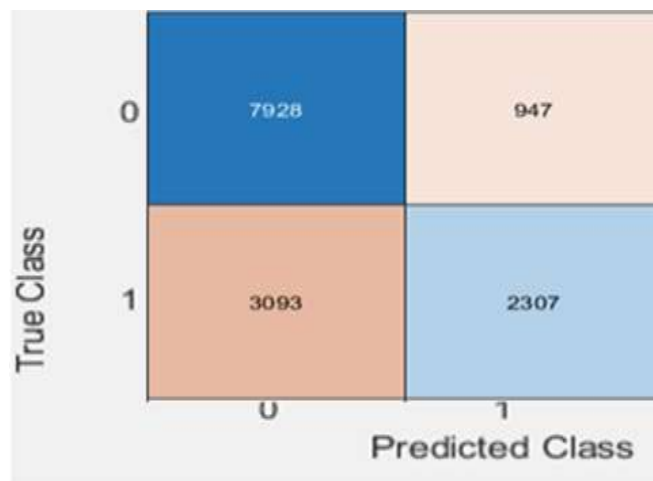


Fig 13: Subspace discriminant

ROC curve of various classifiers:

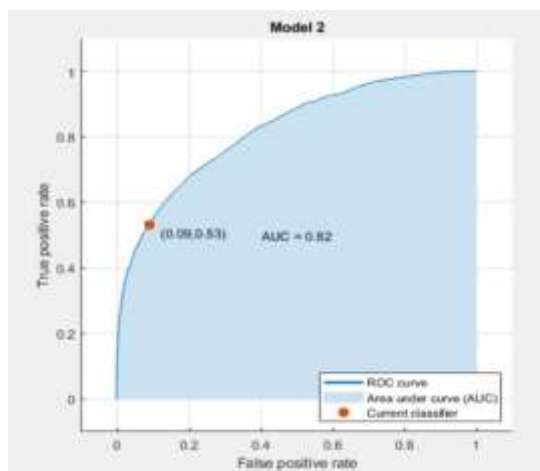


Fig 14: Quadratic SVM (Bangladeshi Bangla)
Bengal Bangla)

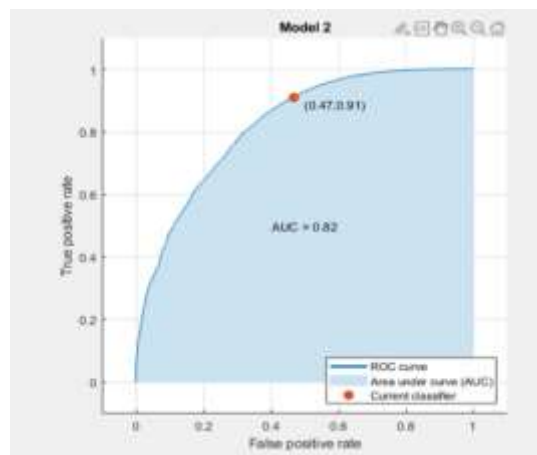


Fig 15: Quadratic SVM (West Bengal Bangla)

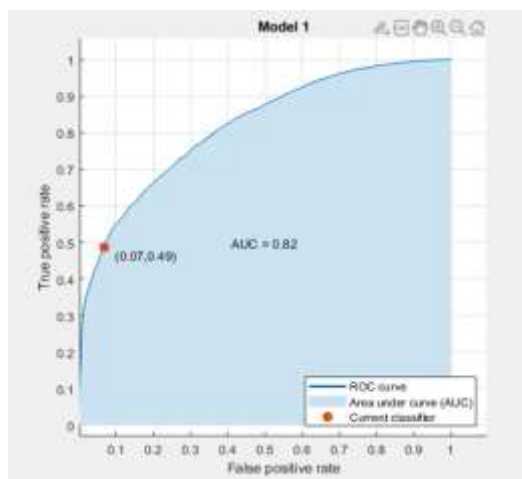


Fig 16: Linear SVM (Bangladeshi Bangla)
Bengal Bangla)

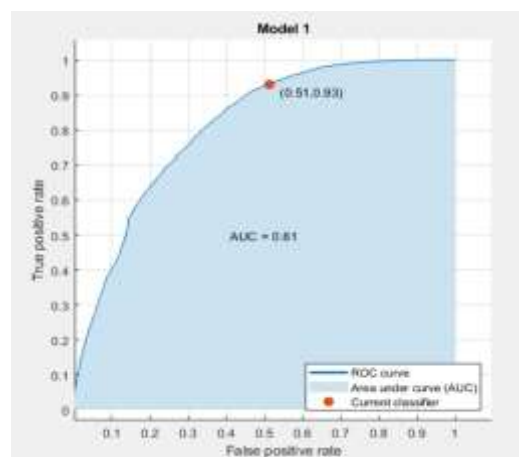


Fig 17: Linear SVM (West Bengal Bangla)

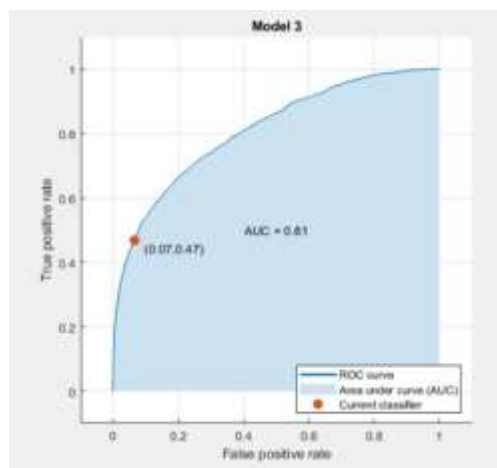


Fig 18: Medium Gaussian SVM
(Bangladeshi Bangla)

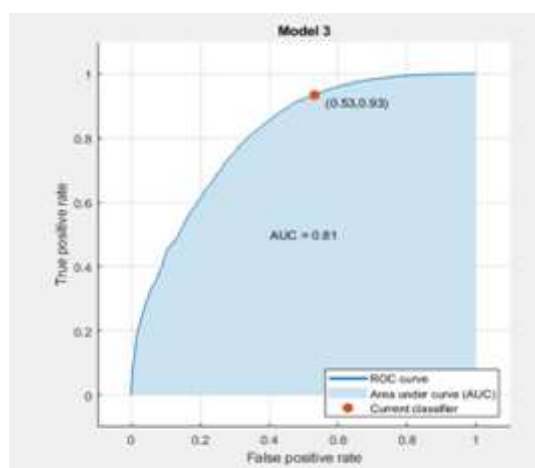


Fig 19: Medium Gaussian SVM (West Bengal Bangla)

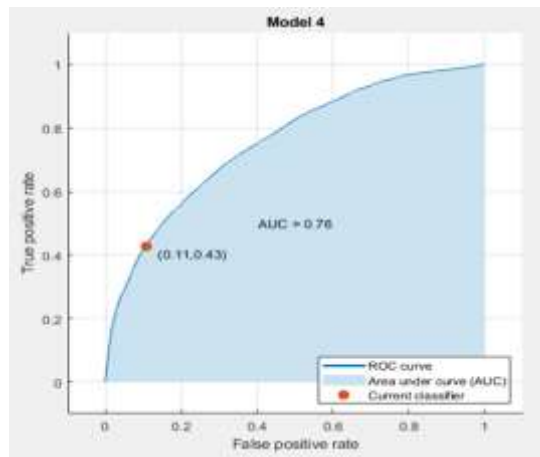


Fig 20: Boosted Trees (Bangladeshi Bangla)

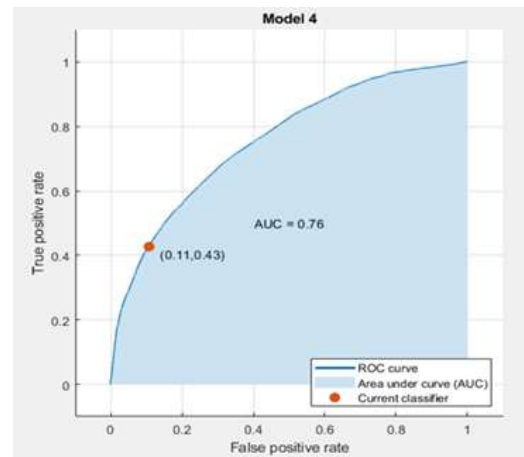


Fig 21: Boosted Trees (West Bengal Bangla)

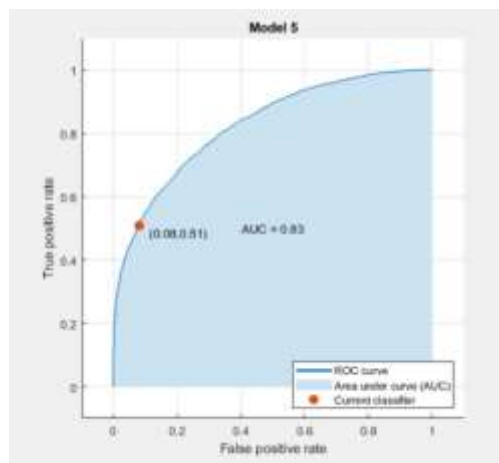


Fig 22: Subspace discriminant (Bangladeshi Bangla)

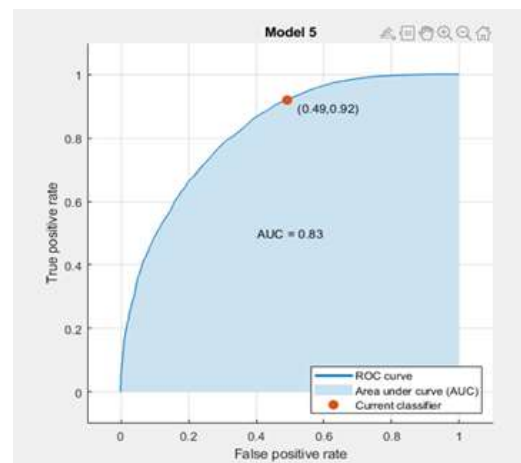


Fig 23: Subspace discriminant (West Bengal Bangla)

In case of ROC curves output is directly related to steepness and in linear SVM we get much steeper curve value. Thus Linear SVM depicts the best functionality among categorizing Bangladeshi Bangla and West Bengal Bangla classification.

CONCLUSION

As we consider whole algorithm, the precise results were given by Linear SVM and gives the expected outcomes. This categorizer assists in classifying languages like Bangla of Bangladesh and Bangla of Western Bengal. This categorizer would prove useful in classifying other accents too and the differentiation of languages will become much easier and understandable.

REFERENCES

- [1] M. S. Erbaugh, "Classifiers are for specification: Complementary Functions for Sortal and General Classifiers in Cantonese and Mandarin," *Cahiers de Linguistique Asie Orientale*, vol. 31, no. 1, pp. 36-69, 2002.
- [2] S. Y. Killingley, *Cantonese classifiers: Syntax and semantics*, Newcastle upon Tyne: Grevatt & Grevatt, 1983.
- [3] N. V. p. S. Sharma, "English to Hindi Statistical Machine Translation System," TIET Digital Repository, 2 August 2011.
- [4] G. C. M. Petrovska-Delacrétaz, "Data Driven Approaches to Speech and Language Processing," in Springer, Heidelberg, 2004.
- [5] D. Roe, F. Pereira, R. Sproat, M. Riley, P. Moreno and A. Macarron, "Efficient grammar processing for a spoken language translation system," in [Proceedings] ICASSP-92: 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing, San Francisco, CA, USA, USA, 1992.
- [6] H. Kitano, "Phi DM-Dialog: an experimental speech-to-speech dialog translation system," vol. 24, no. 6, pp. 36-50, 1991.
- [7] A. Terekhov, "Automating language conversion: a case study (an extended abstract)," in *Proceedings IEEE International Conference on Software Maintenance. ICSM 2001*, Florence, Italy, Italy, 2001.
- [8] J. a. J. K. Halpern, ""Pitfalls and Complexities of Chinese to Chinese Conversion."," in *International Unicode Conference*, Boston, 1999.
- [9] M. S. H. Nuzhat Atiqua Nafis, "Speech to Text Conversion in Real-time," *International journal of innovation and scientific research*, vol. 17, pp. 271-277, 2015.
- [10] H. A. K. S. a. I. Z. Bakr, "A hybrid approach for converting written Egyptian colloquial dialect into diacritized Arabic."," *international conference on informatics and system*, 2008.
- [11] A. l. o. o. p. B. J. W. C. D.RileyAlejandroMacarrón, "A spoken language translator for restricted-domain context-free languages," *Elsevier B.V.*, vol. 11, no. 2-3, pp. 311-319, 1992.
- [12] P. S. Rathod, "Script to speech conversion for Hindi language by using artificial neural network," in *2011 Nirma University International Conference on Engineering*, Ahmedabad, Gujarat, India, 2011.
- [13] B. GRAY, "Grammatical change in the noun phrase: the influence of written language use," *Cambridge University Press*, vol. 15, no. 2, pp. 223-250, 2011.
- [14] K. a. R. K. Kamble, "A review: translation of text to speech conversion for Hindi language."," *International Journal of Science and Research (IJSR)*, vol. 3, 2014.

- [15] N. a. K. A. Swetha, "Text-to-speech conversion," International Journal of Advanced Trends in Computer Science and Engineering , vol. 2, no. 6, pp. 269-278, 2013.
- [16] A. W. Hermann Hild, "Integrating Spelling Into Spoken Dialogue Recognition," in European Conference on Speech Communication and Technology, Germany Carnegie Mellon University, Pittsburgh, USA, 1995.
- [17] B. Sarkar, K. Datta, C. D. Datta, D. Sarkar, S. J. Dutta, I. D. Roy, A. Paul, J. U. Molla and A. Paul, "A Translator for Bangla Text to Sign Language," in 2009 Annual IEEE India Conference, Gujarat, India, 2009.
- [18] A. N. 1. a. J. Z. Merijn Beeksma 1, "shotgun: converting words into triplets: A hybrid approach to grapheme-phoneme conversion in Dutch," John Benjamins , vol. 19, no. 2, pp. 157-188, 2016.
- [19] P. Khilari, "A REVIEW ON SPEECH TO TEXT CONVERSION METHODS," Computer Science, 2015.
- [20] B. temple, "Qualitative Research and Translation Dilemmas," Sage Journals, vol. 4, no. 2, 2004.
- [21] P. S. B. K. Nakib Aman Turzo, "Interpretation of Sadhu into Cholit Bhasha by Cataloguing and Translation System," International Journal Of Trend in Scientific Research and Development, vol. 4, no. 3, pp. 1123-1130, 2020.