

Native Language Identification from English Writing

SNLP TERM PROJECT

UNDER THE GUIDANCE OF PROF. PAWAN GOYAL

GROUP 8

J RATAN RAHUL 10CS30022

P M ARPITA 10CS30041

CH ARCHITHA 10CS30015

YETESH CHOUDHARY 10CS30044

ASHISH VASAVA 10CS30009

PUTTY VIKAS 09CS1044

ARUNENDHER SINGHH 10CS30008

ASHWANI ATTRI 10CS30010

KONTHAM RAJESH 10CS30026

KAUSHIK KUMAR MAHATO 10CS10018

KANWAR PAL DHANDE 10CS10017

SHREYAS KHAIRKAR 11CS30034

Outline

- ▶ Objectives
- ▶ Introduction
- ▶ Dataset Construction
- ▶ Approach Taken
- ▶ Experimental Results
- ▶ Examples
- ▶ Reference

Objective

- ▶ Aim to automatically identify the native language of a writer from its English writings (articles/blogs/essays etc.)
- ▶ Applications:
 - Authorship Profiling
 - Education : more targeted feedback to language learners

Introduction

- ▶ For many years it has been presumed that the only major source for syntactic errors in adult second language performance was the performer's first language.
 - *Today in this article I would discuss about the Database Mail which is used to send the Email using SQL Server.*
 - *In today's article I will discuss Database Mail, which is used to send email using SQL Server.*
 - The first sentence can be very easily translated to an Indian language while the second cannot be translated as easily

Introduction

- Subsequent empirical studies showed that many errors are not traceable to the performer's first language but are common to second language performers of different linguistic backgrounds
 - For example, Nouns in Slavonic languages like Czech and Slovakian do not distinguish between singular and plural.
- These type of errors are found to occur more often than syntactic errors.
- This explains the necessity of supervised learning over reliable corpus dataset for Native Language Identification

Dataset Construction

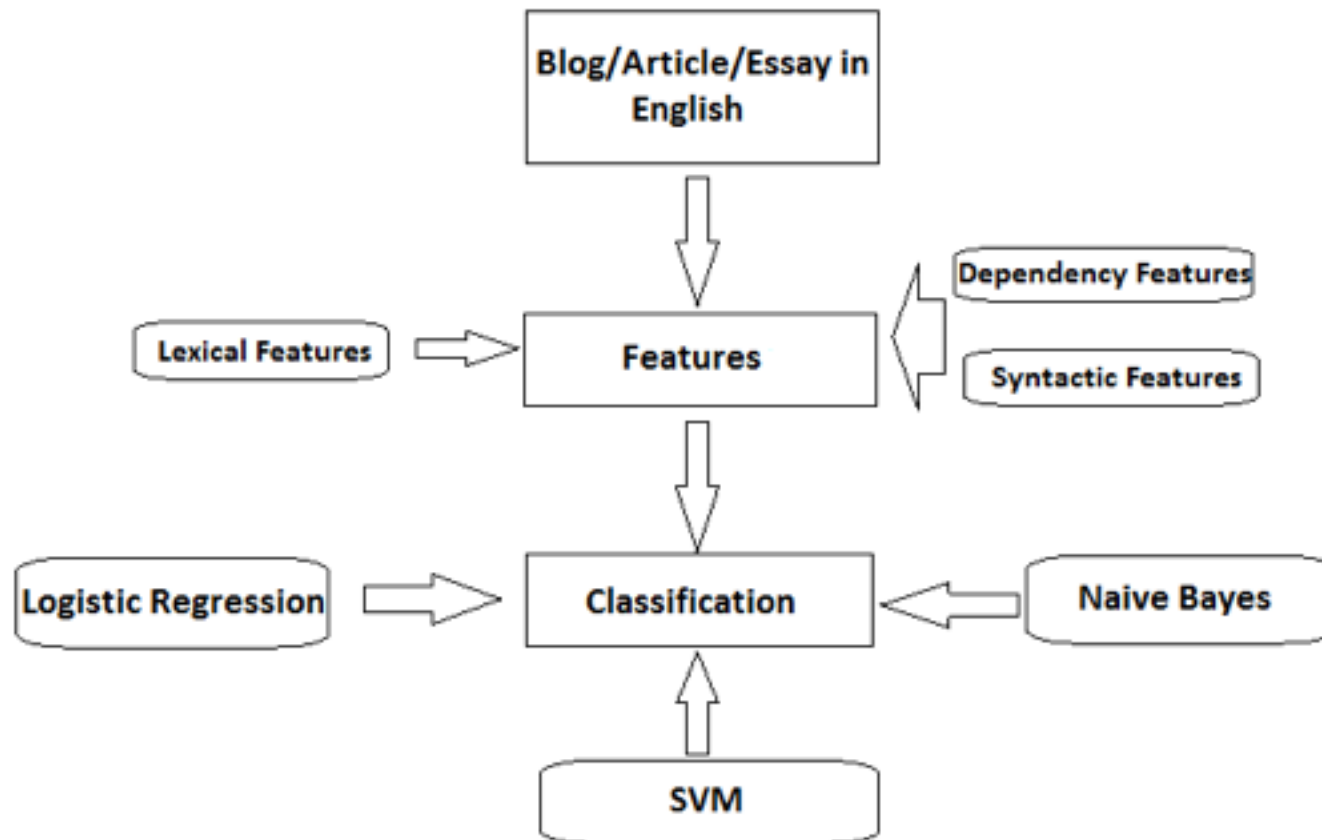
- ▶ Two datasets were used
 - a. Cambridge Learners Corpus First Certificate in English (CLC-FCE)
 - b. The International Corpus Network of Asian Learners of English (ICNALE)
- ▶ In total, 7765 essays from 19 languages were used with training data to test data ratio of 70:30
- ▶ More than 1000 journal entries were extracted from the popular language learner website, lang-8 for using as test data.

Dataset construction

▶ Languages include:

- Catalan
- Chinese
- English
- French
- Filipino
- German
- Greek
- Indonesian
- Italian
- Japanese
- Korean
- Polish
- Portuguese
- Russian
- Spanish
- Swedish
- Thai
- Turkish
- Urdu

Flow Diagram



Approach Taken

- ▶ The problem is addressed as a supervised multi classification task.
- ▶ Trained our data on feature set using 4 classification models.
 - ▶ Logistic Regression
 - ▶ Gaussian Naïve Bayes
 - ▶ One vs rest classifier
 - ▶ Support Vector Machine
 - ▶ Kernel – 'rbf', tol = '0.001'
 - ▶ Multinomial Naïve Bayes
 - ▶ One vs rest classifier

Approach Taken

- ▶ The feature set comprises of the following features:
 - ▶ **Functional Words:** Certain functional words are more common in one language compared to the others. The functional words include The, to, I, and, a, was etc. 45 functional words are used as features.
 - ▶ **Word n-gram:** Writer's native language influences their choice of words. So, 1-gram and 2-gram words are used for calculating the features. 70 features are extracted from 1-gram and 290 features from 2-gram.
 - ▶ **Use of Punctuation:** The speakers of different languages use punctuation in different ways. The following 2 features are considered:
 - ▶ The number of punctuation marks used for sentence.
 - ▶ The number of punctuation marks used for word.

Approach Taken

- ▶ **Number of Unique Stems:** Speakers of different native languages differ in the amount of vocabulary used. The relative frequency of the total number of unique stems is considered as a feature.
- ▶ **Misuse of articles:** The number of instances in which an article is inconsistent with the plural and uncountable nouns are considered as features.
- ▶ **Missing Punctuation:** The relative frequency of missing punctuation after introductory (however, furthermore) and subordinating conjunction (after, before, even though) phrases are considered as features.

Approach Taken

- ▶ **Words per sentence:** Number of words per sentence is considered as a feature.
- ▶ **Tense and Aspect Frequency:** Three tenses (present, past, future) are considered for calculating tense frequency. Four aspects (simple, perfect, progressive, perfect progressive) are considered for calculating aspect frequency. These frequencies are considered as features.
- ▶ **Part-of-speech:** Bigram and trigram pos tags are used for calculating features. Total of 101 bigram pos tags and 131 trigram pos tags are extracted as features.
- ▶ **Passive Constructions:** The count of number of times an author uses passive constructions (count of nsubjpass) is considered as a feature.

Results

- ▶ Used a supervised multi class classifier approach for modeling the data
- ▶ Ratio of Training: Test set = 70:30
- ▶ Total of 5462 training features and 2303 test features sampled over all 19 languages.

Method	Training Accuracy	Test Accuracy
<i>Logistic Regression</i>	100 %	55.97 %
<i>Gaussian Naïve Bayes</i>	96.33 %	52.36 %
<i>SVM</i>	95.97 %	23.75 %
<i>Multinomial Naïve Bayes</i>	59.08 %	13.28 %

Reference

- ▶ Abu-Jbara, Rahul Jha, Eric Morley and Dragomir Radev, Experimental results on the Native Language Identification Shared Task, *Proceedings of the Eighth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 82–88, 2013.