

Similar Language Detection - Report

Group Members:

Naveen Bhookya (140050034),

Chanukya Vardhan (140050043),

Uday Kusupati (140050048),

Raghuveer Uppara (140050050).



Overview :

Discriminating between similar languages is the whole purpose of the project.

We have used a hierarchical method, we first divide all the languages into few groups and next we label the languages in each group.

The first step is Word-Frequency method with tf-idf-qf weighting and second method uses SVM's.

Problem Statement and Description :

The training data consists of sentences labeled with 13 languages which are Bulgarian (bg), Macedonian (mk), Bosnian (bs), Croatian (hr), Serbian (sr), Czech (cz), Slovak (sk), Spanish (esar), Peninsular Spanish (eses), Brazilian Portuguese (ptbr), European Portuguese (ptpt), Indonesian (id), Malay (my) and a unknown language(xx),

Each language has 18,000 sentences for training set and 1,000 sentences for testing. Labelling a query sentence with correct language is the problem.

We divide the training set for each language into 16,000 and 2,000 sentences for training and validation data respectively

Motivation:

Language detection has many uses in various domains like translation, language interpretation etc.

Discriminating between disparate languages is trivial in the sense that it has well known methods for solving.

But the same methods used for disparate languages doesn't yield good results for discriminating between similar languages and this project address this problem

Approaches:

Classification into groups

From the training data, for each language we first parse through all the sentences and find the frequencies of all the words present.

Among all the words and their corresponding frequencies generated for each language, we choose the top N words based on their frequency.

Define x_l for each language l as vector $\in \mathbb{R}^N$, x_l = 1000-tuple of frequency of the chosen words.

Now given an input statement s we define $x(s)$ for language l as follows $x_l(s) \in \mathbb{R}^N$ first initialised to all zeros

$x_{l,i}(s) = 1$ if the i th word of the l th language is present in the input sentence.

Label for statement s is $\text{label}(s) = \underset{l \in \{1 \dots 13\}}{\text{argmax}} (\text{dotproduct of } x_l(s), x_l)$.

The intuition is clear, we choose top words from each language and if a top word is found in the query sentence, we add its frequency to the score and the language with maximum score is label.

From the confusion matrix after first stage, it is clear it is not appropriate since it needs lot of groups to make group accuracy high and 'pt-Br' and 'eses' cannot be put to same group as 'eses' does not contain any 'pt-Br' while 'pt-Br' has lot of 'eses'

Innovation

Solution we thought of is to use tf-idf scoring i.e. the x_i for each language has $\log(1+1/f_l) * \log(1+f_t)$ and $x_i(s)$ has $f(t,q) * \log(1+1/f_l)$ where f_l = fraction of languages which contains the word and f_t = frequency of the word in the language and $f(t,q)$ is frequency word in query ($f(t,q)$ is normalised). The above idea worked and the intuition is the term $\log(1+1/f_l)$ decrease the weight for words occurring in all languages. We choose $N = 5000$ clear from graphs and confusion matrix and group accuracy is greater than 99 percent

Results

The languages bg, cz, mk, sk are already classified with accuracy greater than 99 percent and the other languages has less accuracy between 70-90 percent and we group them as follows bs,hr,sr esar,eses id,my ptbr,ptbt.

Note that languages are classified into corresponding groups with accuracy close to 100 percent which is we intended.

Now each group classification is a different problem and we intend to use ensemble methods with SVMs. For each group we have ensemble of 5 SVMs

Classification within groups:

Our basic idea for feature modeling for the SVM was through word and character n-grams. Character 1,2-grams when enumerated are very less, so they don't make good features. We chose character n-grams with $n \in 3,4,5,6,7,\dots$ and word n-grams with $n \in 1,2,3,\dots$. We needed to model an efficient feature vector from the n-grams with an appropriate dimension.

Since character and word n-grams capture different properties, we felt instead of clubbing them in the same feature vector of an SVM, training an ensemble of SVMs over the pure n-gram features is better.

SVMs for different n-grams as listed above were trained and the overall ensemble accuracy was noted for incremental subsets of 3,4,5,6,7... for characters and 1,2,3,... for words. From the plots we observe that the optimal features would be character 3,4,5-grams and word 1,2-grams, a partial reason also being overfitting for character n-grams with $n > 4$

Feature Engineering (Innovation):

Even after choosing the features as above, the number of features are very high and need to be reduced. For capturing the important features, we score the features using tf-idf scoring with

$$\text{tfidf}(t, D) = \text{tf}(t, D) \times \log(N/\text{df}(t, D))$$

We rank the features according to the tf-idf scores and pick the top n features. n is chosen by plotting the validation accuracies. The accuracies get saturated for $n > 20000$

Given a query, the feature vector is a boolean vector that captures the presence of the n-ranked features. We represent the feature vector in a sparse way as a dictionary

Conclusions:

Test accuracy is calculated for 1000 test examples for each language completely not known to the trained model

The validation and test accuracies are in consensus as in the above table

Although the languages bg,mk can be grouped and similarly cz,sk we didn't group them as the initial classification had very high validation accuracy.

For the group id, my after the initial classification itself the validation accuracy for each language was high as seen in the confusion matrix. But we grouped them for a better accuracy and obtained an increase

Since the grouping accuracy was very close to 100% we now compare the validation accuracies of both the stages for each language as shown below

Results:

The results for the language classes are as follows

Language group	Validation Accuracy(%)	Test Accuracy(%)
bs, hr ,sr	88.48	89.6
eses, esar	89.55	90.1
id, my	98.82	98.6
ptbr, ptbt	91.72	92.4
bg	99.4	99.6
cz	99.4	99.3
mk	99	99.8
sk	99.5	99.6

Comparison:

Language	Stage-1(%)	Stage-2(%)
bs	70	84.5
hr	83	88.8
sr	90	95.6
eses	90	88.8
esar	75	91.7
id	93	98.5
my	96	98.7
ptbr	87	91.1
ptbt	87	94.1
bg	99.4	99.4
cz	99.4	99.4
mk	99	99
sk	99.5	99.5

Future Work and Improvements:

The accuracy of eses has fallen due to grouping with esar but the grouping is necessary so as to increase the accuracy of esar which actually increased drastically. More state of art ensemble methods and classifiers can be tested and also the intial grouping can be done faster with smaller bag of words by better tf-idf-qf scoring which we couldn't do due to the limited scope of this project

We need to appreciate the fact that the accuracy of bs and esar have drastically increased due to grouping. Though the accuracy of bs increased very well, the overall accuracy for the language group bs, hr, sr is still low and it should be improved by different techniques

Contributions:

Raghuveer U - Coding for Language grouping, SVM modeling

Chanukya G - Coding for Language grouping, Feature Engineering for SVM

Uday K - Implementation Design and Analysis in both stages, Usage of SVM

Naveen B - Tuning parameters from validation accuracies

References:

We will be working on the data set that was provided in the competition. There are 18000 training samples for each language

<https://github.com/Simdiva/DSL-Task/tree/master/data/DSLCC-v2.0>.

Reference : http://cs229.stanford.edu/proj2015/335_report.pdf