

# STATISTICAL METHODS IN ARTIFICIAL INTELLIGENCE

Team 26 - BOLTS

---

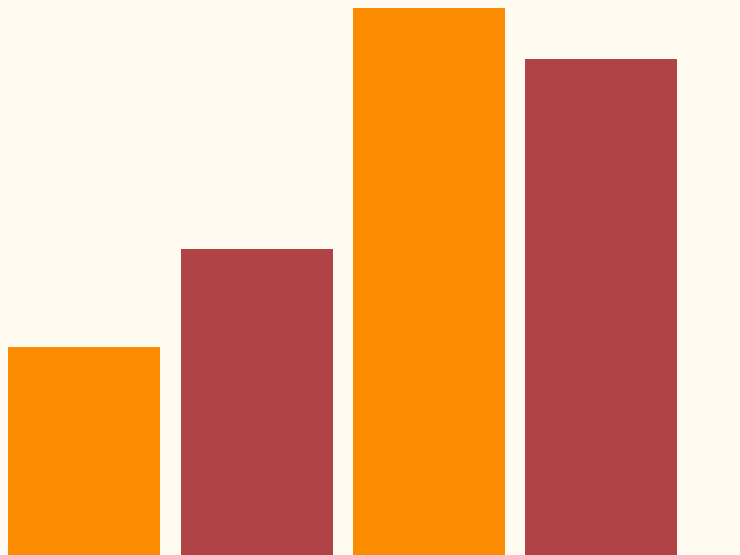
Nishanth Reddy

Sravya Varma

Poojitha Nandigam

# THE PROBLEM

Discriminating Similar Language (DSL), to predict the language of sentences written in similar languages



# DEFINITION

The discrimination of similar languages can be defined as the subtask of the Language identification problem. Language Identification is a fundamental task in the area of natural language processing.

Unlike well-separated languages, multilingualism, dialects of language can seriously degrade the quality of Language Identification. Discrimination of Similar Languages, noisy data, non well-formatted text, short sentences, mixed language are other examples of challenging problems in this field.

.....

# DATASET

The data sets we used were part of the DSL-Shared Task of 2015. The task provided participants with training (17000 examples per language), development (2000 examples per language) and testing sets (1000 examples per languages). The sets consisted of individually labeled sentences extracted from the journalistic corpora in 13 different languages. The languages were part of 6 language groups.

# Method

## Step1

Determine language group  
by word frequency method

## Step 2

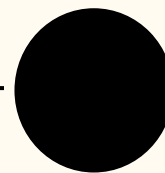
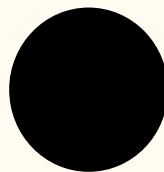
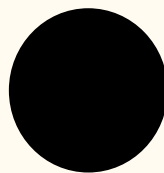
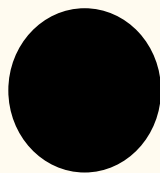
Train group-specific SVMs  
with character and word  
n-grams using training  
dataset.

## Step 3

Rank features with tf-idf  
scoring and tune  
parameters using devel  
dataset

## Step 4

Combine decisions to  
predict test dataset using  
ensemble methods to  
reduce variance

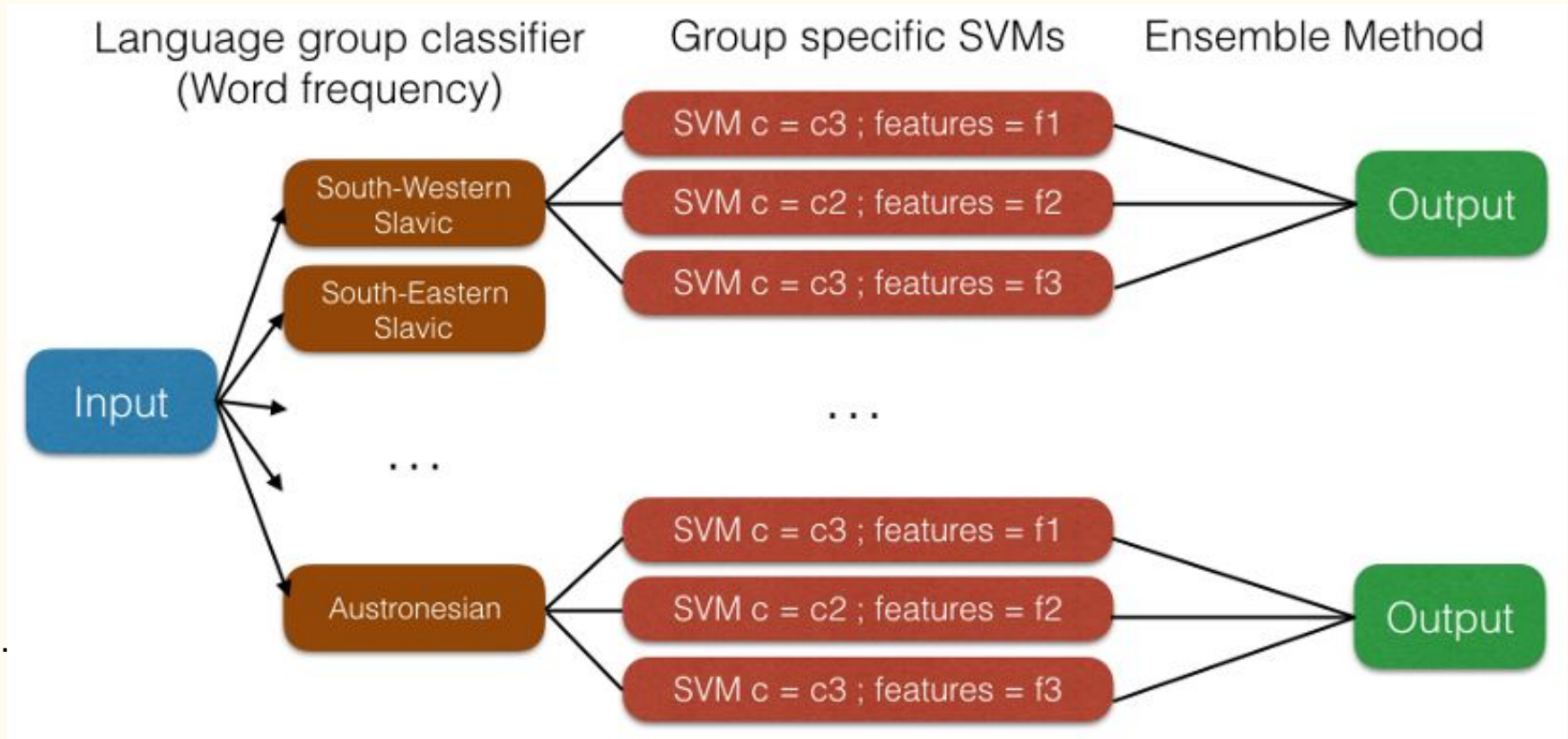


An aerial photograph of the New York City skyline at dusk. The sky is a mix of dark blue and orange, with scattered clouds. The city is densely packed with skyscrapers, many of which are illuminated with lights. The Empire State Building is prominent in the center, with its top lit in red and green. The Hudson River is visible on the right side of the image.

# LIBRARIES

- LIBLINEAR SVM to implement multi-class SVM

# THE SOLUTION



# SCOPE OF THE PROJECT

- DSL Shared Task of 2015 made available the required corpus.
- Now, word frequency method is applied to group the data into 6.
- To differentiate between the similar languages present in the group , Multiclass SVM is used .
- To avoid curse of dimensionality, tf-idf scoring is used for feature scoring, and only high ranked features are considered for ensemble.
- Ensemble combines various SVMs with different hyperparameters. Ensemble method used is Mean Confidence.

.....



# INTERIM EVAL II

- Implementation of
  - Word - Frequency
  - Multi class -SVM
  - tf-idf
  - Ensemble (Tentative)

# FINAL SUBMISSION

- Implementation of the above mentioned method and 2nd method (From other paper);

THANK YOU

