

# Similar Language Detection

## Group Members:

Naveen Bhookya(140050034), Chanukya Vardhan(140050043),  
Uday Kusupati(140050048), Raghuveer Uppara(140050050).

## Stage 1 Report:

- From the training data, for each language we first parse through the sentences and find the frequencies of all the words appearing in all the sentences.
- Among all the words and their corresponding frequencies generated for each language, we choose the top 1000 words based on their frequency values.
- Define  $x^l$  for each language  $l$  as vector  $\in \mathbb{R}^{1000}$ ,  $x^l = 1000$ -tuple of frequency of the chosen words.
- Now given an input statement  $s$  we define  $x^l(s)$  for language  $l$  as follows:  
 $x^l(s) \in \mathbb{R}^{1000}$  first initialised to all zeros.  
 $x_i^l(s) = 1$  if the  $i$ th word of the  $l$ th language is present in the input sentence.
- Label for statement  $s$  is  
$$\text{label}(s) = \underset{l \in \{1,2,\dots,13\}}{\text{argmax}} (\text{dotproduct of } x^l(s), x^l).$$
  
The intuition is clear.

## Observations:

- The confusion Matrix generated after running the code on training data(18000 samples for each language) is as follows:

$C =$	—	<i>bg</i>	<i>bs</i>	<i>cz</i>	<i>esar</i>	<i>eses</i>	<i>hr</i>	<i>id</i>	<i>mk</i>	<i>my</i>	<i>ptbr</i>	<i>ptbt</i>	<i>sk</i>	<i>sr</i>
	<i>bg</i>	17379	0	0	0	0	0	0	620	0	0	0	0	1
	<i>bs</i>	0	1232	34	0	48	288	0	0	0	51	22	25	16299
	<i>cz</i>	0	396	12486	0	199	128	0	0	1	1115	411	1843	1421
	<i>esar</i>	1	0	0	0	17999	0	0	0	0	0	0	0	0
	<i>eses</i>	2	0	0	0	17998	0	0	0	0	0	0	0	0
	<i>hr</i>	0	1872	57	0	73	1146	0	0	0	92	31	35	14694
	<i>id</i>	2	0	1	0	29	0	14879	0	3077	7	4	0	1
	<i>mk</i>	36	1	1	0	0	0	0	17961	0	0	0	0	0
	<i>my</i>	0	0	9	0	25	0	2227	0	15715	1	1	2	2
	<i>ptbr</i>	1	0	0	0	9819	0	0	0	0	7557	622	0	1
	<i>ptbt</i>	1	0	0	0	10879	0	0	0	0	5822	1297	0	1
	<i>sk</i>	0	75	151	0	62	49	0	0	0	943	349	15769	602
	<i>sr</i>	0	214	5	0	59	38	0	0	0	38	12	29	17605

- This will be used for dividing the languages into groups.