# Assignment II & III: POS Tagging

Rajesh Kumar (201307556)

Kiran Raj S R (201307593)

Abhinav Yadav (201307635)

October 22, 2013

This project consists of 2 parts:

1. POS Tagging using supervised learning.
2. POS Tagging using unsupervised learning.

We have selected *Hindi, Telugu, Kannada & Tamil* languages for our project.

# 1 POS Tagging using Supervised Learning

## 1.1 Assumptions

### (i) Tag Set

We have used LTRC tag set which consists of 26 tags.

### (ii) Start and End Markers

Each sentence has been enclosed within <s> START and </s> END markers. The role of these tags is just to avoid inclusion of invalid transitions between the sentences. And hence they don't appear in the output file.

**Example:**

<s> START आपके PRP हिन्दी NN पसन्द NN करने VM पर PSP खुशी NN

**हुई VM** </s> END

<s> START **साथ NS** Tही RP मैने PRP लेख NN पर PSP कुछ QF

बदलाव NN किए VM है VAUX </s> END

Without START and END tags, VM to NS transition (shown above in bold) could have been a valid transition, which has been ignored now.

### (iii) Handling unknown words

Whenever an unknown word is encountered in test data, its emission value is 0. Due to which 0 gets propagated while creating probability matrix in Viterbi algorithm. We tried stemming for this earlier, but they don't seem to give an appreciable result w.r.t unknown words. So to avoid this problem, we have used *emission value=0.001* for all the unknown words encountered in testing sentences.

## 1.2  Task Descriptions

### (i)  Manual tagging

We performed manual tagging on the training corpus provided to us to get the POS tags corresponding to each word. This tagging was performed with the help of shallow parser provided by LTRC.
We cleansed the output of shallow parser to obtain the tags (by simple script) and manually cross verified the tags assigned by shallow parser.

### (ii)  Creating Emission Matrix

Using the data provided by the files created in step-1, we can create emission matrix as follows.

$$P(w_i | g_i) = \text{count}(w_i, g_i) / \text{count}(g_i)$$

### (iii)  Creating Transmission Matrix

Using the data provided by the files created in step-1, we can create transmission matrix as follows.

$$P(g_i | g_{i-1}) = \text{count}(g_i, g_{i-1}) / \text{count}(g_{i-1})$$

### (iv)  Creating Viterbi Matrix

We create a 2 dimensional matrix of **k*n** size.
k: Number of tags.
n: Number of words in the test sentence under consideration.

In this matrix, rows correspond to the tags and the columns correspond to the words in the test sentence under consideration.

We use the following method to fill the values in this matrix.

$$P(g_i/w_i) = P(w_i/g_i = t^1) * \mathbf{max}(\, P(g_i = t^1/g_{i-1} = t^k) * P(w_{i-1} / g_{i-1} = t^k)\,)$$
$$\text{Where } k = 1, 2, 3 \ldots \text{number of tags}$$

**(v)    Traversing the Viterbi Matrix**

We traverse the matrix generated in the previous step to get the tags corresponding to each word for the test sentence under consideration.

## 1.3    Results

After applying the steps mentioned above, we obtained the pos tags corresponding to each word for all the test sentences (for 4 languages).

The corresponding output files can be found at,
./saved outputs/Supervised Tagging/<language>_ tags.txt

Sample output is as follows.

**Hindi:** हिंदी_NN के_PSP पहले_NST प्रगतिशील_JJ लेखक_NN थे_VM प्रेमचंद_NN ._.

**Telugu:** హిందూ_NNP సంఘం_NN ఒక_JJ కులాల_NN కూటమి_NN ._.

**Kannada:** ಭಕ್ತ_NN ಕನಕದಾಸ_NNP ಚಿತ್ರದಲ್ಲಿ_NN ಸಾಹಿತ್ಯದಲ್ಲಿ_NN ಸಹಾಯಕರಾಗಿ_NN ದುಡಿದರು_NN ._.

**Tamil:** தலித்_NN இலக்கியம்_NN பற்றி_RB பல_QF கட்டுரை_NN எழுதியுள்ளார்_VM ._.

## 1.4  Observations

### (i)    Execution time

**Table:** Training time for different number of training sentences.

| Number of Sentences | Time taken for training (in Seconds) |
| --- | --- |
| 200 | 0.1170 |
| 400 | 0.2899 |
| 600 | 0.4769 |
| 800 | 0.7309 |
| 1000 | 0.9470 |

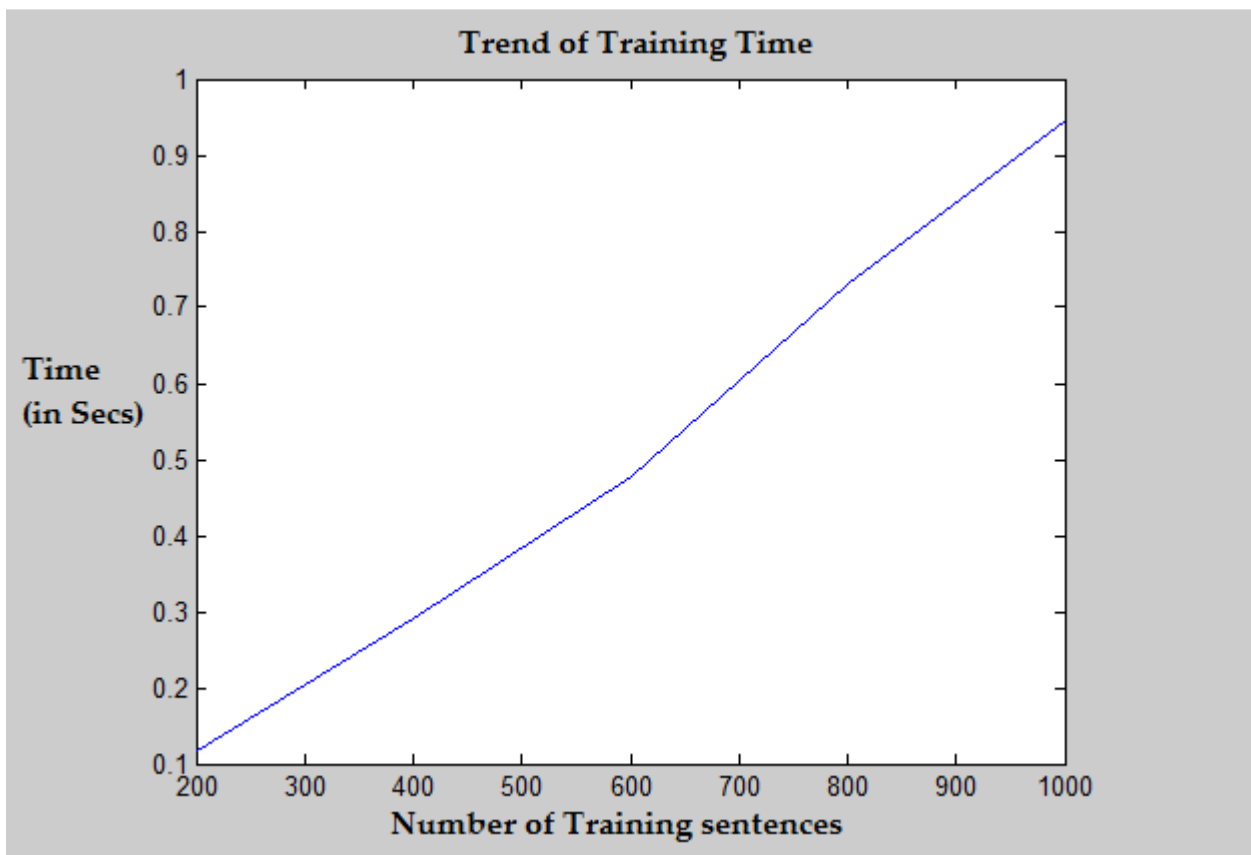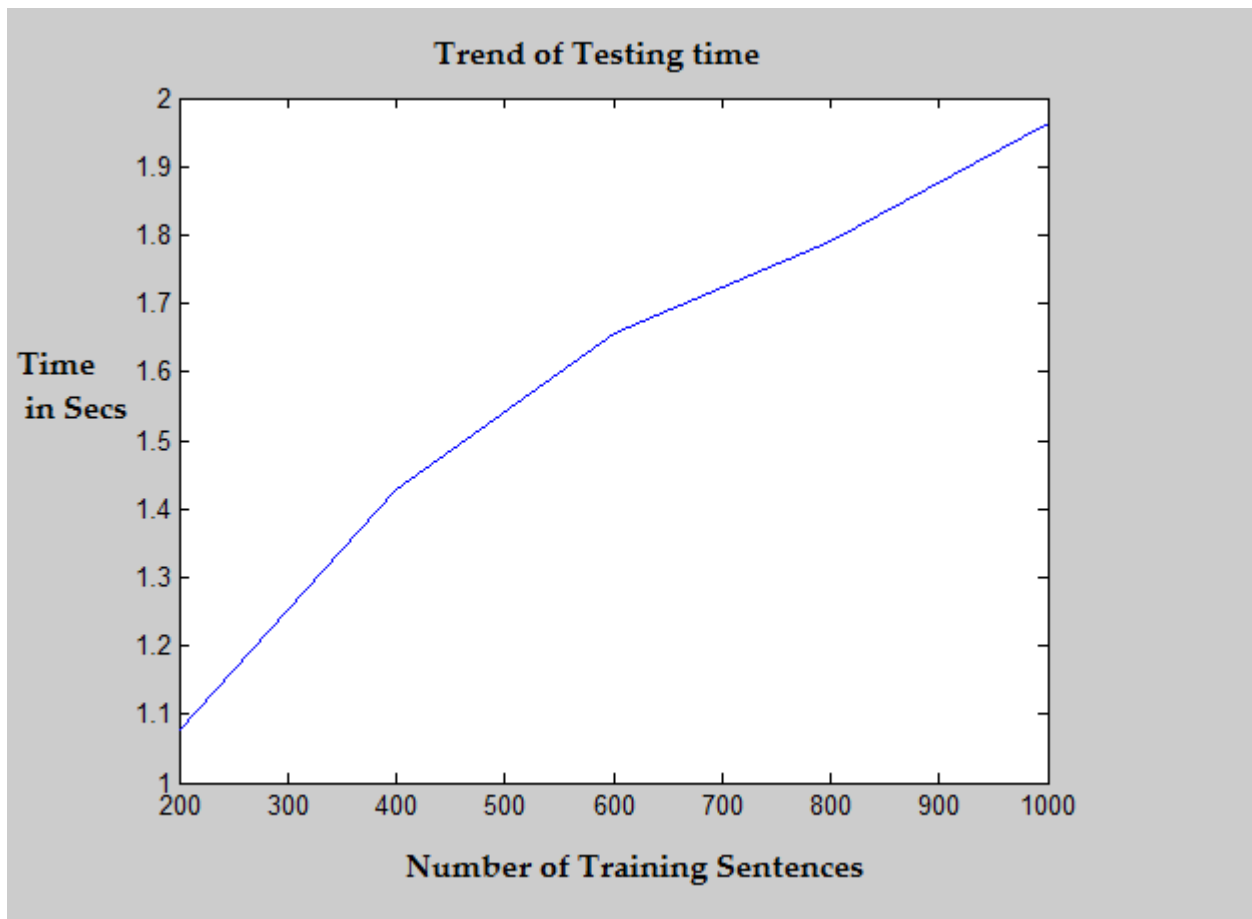**Figure:** Training time for different number of training sentences.

**Table:** Testing time for different number of training sentences.

| Number of Training Sentences | Time taken for testing 100 Sentences (in Seconds) |
|---|---|
| 200 | 1.077 |
| 400 | 1.4270 |
| 600 | 1.6559 |
| 800 | 1.7909 |
| 1000 | 1.9649 |

**Figure:** Testing time for different number of training sentences.



Average time taken to tag a single test sentence is 21.19 milliseconds.

### (ii)   Stemming

To handle unknown words, we tried to stem a word to make it a length n word [n=4, 6]. But number of unknown words didn't drop significantly for the given corpus. Rather it introduced a few problems as explained below.

Example: For n=4, "Import" and "Important" will give us the same stem "Impo", which is undesirable since these words have completely different meaning.

So a blunt stemming is not effective. A sensitive stemming (with linguistic rules) might have given a better result. But keeping in mind the complex linguistic properties of words in unknown languages, we preferred **NOT** to use Stemming for our project.

# 2    POS Tagging using Unsupervised Learning

## 2.1   Assumptions

### (i)   Number of Clusters

Number of clusters considered by us in this case =26 i.e. Number of tags in the LTRC tag set. [i.e. 1 cluster corresponding to each tag]

### (ii)   Context Depth

For checking the context, we have considered 4 words to the left of the given word (left context) and  4 words to the right of the given word (right context).

### (iii)   Initialization of vectors

To introduce randomness in the initialization of 26 centroids (while using K-means clustering), we have initialized the centroids with feature vectors of random words (picked from our training corpus).

### (iv)   Number of feature words

We have used top 125 words as feature words. But we have removed top 25 words (Stop Words) among these 125 words. So only the next 100 words have been considered as feature words.

### (v)   Distance measure

We have used euclidian distance for calculating the distance between two feature vectors.

### (vi)   K-Means Clustering methodology used

## 2.2 Task Description

**(i)    Feature word selection**

From the given training corpus, we have considered top 100 words as feature words (after chopping off top 25 words (stop words) from top 125 words).

**(ii)   Creation of feature vectors**

We have generated a *m x 2n* co-occurrence matrix where [m: number of words, n: number of feature words]

a.      $C_{i,j}$ : Count of number of times $i^{th}$ token type occurs with the $j^{th}$ token type as left hand neighborhood of length 4. (left context).

b.      $C_{i,j+n}$ : Count of number of times the same $i^{th}$ feature token occurs with the same $j^{th}$ token type as right hand neighborhood of length 4 (right context).

**(iii)  Comparing words**

Each word will be represented by a feature vector. Distance between each word is given by Euclidian distance.

**(iv)   K-Means Clustering on feature vectors obtained.**

a.      Centroid initialization

b.       Shuffle the centriods on the basis of average distance. Check the Error.

c.       Continue till the error becomes very less.

**(v)    Using HMM to get the tags**

Each word will belong to a cluster. Now we change the tag set used in phase-1 (supervised tagging) of our project, we can get the corresponding tags for testing data. [By application of same HMM model used in phase-1]

## 2.3 Results

Words in same box belong to same cluster.

We are enumerating top 10 words from any 10 non-empty clusters. [Less than 10 if the cluster does not have enough words].

# Hindi:

| | | | |
|---|---|---|---|
| ऋषि | भाषाओं | कवि | पूर्णिमा |
| प्रधान | तमिल | भाषा | सदस्य |
| प्रधानमंत्री | शिक्षा | संस्कृति | अनाम |
| सेनानी | लेखक | प्रथम | अतुल |
| वकील | उर्दू | सम्पादन | पृष्ठ |
| पत्र | अरबी | नाटक | विजय |
| कविता | लिखा | पत्रकार | ठाकुर |
| कहानी | ज्ञान | वाणी | अरविन्दन् |
| नर | लेखो | कोश | युकेश |
| सखा | अंग्रेजी | श्रेष्ठ | वर्मन |
| | | | भटनागर |
| संस्कृत | राज्य | उनके | विश्वविद्यालय |
| उपन्यास | उत्तर | हेतु | ऋग्वेद |
| जैसे | भारतीय | यात्रा | आज |
| प्रयाग | सिक्किम | गाइड | मानते |
| माया | जिला | संख्या | छोटी |
| विचार | प्रदेश | दो | इकाई |
| काऊ | | इतिहास | रचना |
| विद्वान | | करीब | बड़ा |
| लगता | | त्याग | जनसंख्या |
| मूल | | | पुराने |
| कि | सनातन | | |
| प्रकार | वैदिक | | |
| हो | प्रारंभिक | | |
| कर | शैली | | |
| तरी | संसार | | |
| मी | अर्थशास्त्र | | |
| मे | परमात्मा | | |
| ये | मिलन | | |
| कला | भरत | | |
| देखें | पंकज | | |

## Telugu:

| | | | |
|---|---|---|---|
| రవాణా | ఆ | నుంచి | కడప |
| జాబితా | తరువాత | వరకు | చిత్రం |
| పేజీకి | హైదరాబాదుకు | రైలు | ఆస్ట్రేలియా |
| దీని | బదిలీ | నుండి | పక్షి |
| అప్పోడు | అయ్యాడు | బస్సు | ఉన్నవి |
| విజయవాడ | గ్రామంలో | గ్రామానికి | సమాచారం |
| చాలా | ప్రజల | తరచు | టెస్ట్ |
| పాత | ప్రధాన | నుంచి | కడప |
| అని | వృత్తి | వరకు | చిత్రం |
| ఆంధ్రప్రదేశ్ | వ్యవసాయం | రైలు | ఆస్ట్రేలియా |
| | | | పక్షి |
| జిల్లా | కేవలం | వికీపీడియా | ఊరు |
| రాజ్ | ఇది | వ్యాసం | వికీపీడియాలో |
| మంత్రిత్వ | నాకు | పేజీలో | లేదు |
| శాఖ | ఎందుకంటే | సినిమా | మధ్య |
| వెబ్ | ఇంజనీరింగ్ | వారపు | పేజీ |
| సైటులో | ప్రయోగాలు | | మంచినీటి |
| తాలూకాల | వికీలో | | సంవత్సరం |
| జిల్లా | అర్థం | | ఉన్నది |
| రాజ్ | మౌలిక | | ఊరిలో |
| మంత్రిత్వ | ఉంటుంది | | ప్రతి |
| శాఖ | | | |
| ఉదాహరణకు | దేశం | | |
| వికీపీడియా:తొలగింపు | ప్రభుత్వం | | |
| వ్యాసాలు | రాష్ట్రంలోని | | |
| చర్చ | దేశము | | |
| పంట | | | |
| వరలక్ష్మీ | | | |
| వ్రతం | | | |
| అన్న | | | |
| విధమైన | | | |
| భౌతిక | | | |

Kannada:

| | | | |
|---|---|---|---|
| ಸಂಯುಕ್ತ | ಗಳ | ರಾಷ್ಟ್ರ | ಕೆರಿಬ್ಬಿಯನ್ |
| ಸಂಸ್ಥಾನ | ಸಾಂಸ್ಕೃತಿಕ | ವಿಶ್ವವಿದ್ಯಾಲಯವು | ಸಮುದ್ರ |
| ಸಂವಿಧಾನ | ಮೈಸೂರು | ಇದೆ | ಭಾಗದಲ್ಲಿ |
| ದಿನಾಚರಣೆ | ಭಾರತದ | ಜಾನಪದ | ಇರುವ |
| ಖಂಡದ | ಅನೇಕ | ತಜ್ಞ | ದ್ವೀಪ |
| ಅತ್ಯಂತ | ಆದ್ದರಿಂದ | ಅಕಾಡೆಮಿ | ಜನಸಂಖ್ಯೆ |
| ಚಿಕ್ಕ | ಹಾಗೂ | ಪ್ರಶಸ್ತಿ | ಹೊಂದಿರುವ |
| ದೇಶ | ಜಿಲ್ಲಾ | ರಾಜ್ಯೋತ್ಸವ | ಒರೆ |
| ಆಗ್ನೇಯ | ಕೇಂದ್ರ | ಮಹಿಳಾ | ಅಕ್ಷರಗಳು |
| ಅನನ್ಯ | ಪಶ್ಚಿಮ | ಸಂಗೀತ | ಹೀಗೆ |
| | | | ಕೆರಿಬ್ಬಿಯನ್ |
| ಪೂರ್ವ | ಅಮೇರಿಕ | ಅವರ | ಟೆಂಪ್ಲೇಟ್ |
| ಕಡಿಮೆ | ಉತ್ತರ | ಹೆಸರಿನ | ಕೆಲವು |
| ಇದು | ಮಧ್ಯ | ಇದೇ | ಮೇಲೆ |
| ಕೊಡಲಾಗಿದೆ: | ಚಿತ್ರ | ಕಾದಂಬರಿ | ಹಲವು |
| ಮರಗಳು | ದಕ್ಷಿಣ | ಆಧಾರಿತ | ಈಗ |
| ಬಹಳ | | | ಮಾಹಿತಿ |
| ಆದರೆ | | | ಬೆಂಗಳೂರು |
| ಅವರು | | | ಹುಬ್ಬಳ್ಳಿ |
| ತಮ್ಮ | | | ಪತ್ರಿಕೆಗಳಲ್ಲಿ |
| ಮಾಡಿ | | | ವಿಶ್ವಕೋಶದ |
| ಸಂಸ್ಥೆಯ | ಫಿಲಿಪ್ಪೀನ್ಸ್ | | |
| ಸೋದರ | ಪ್ರಮುಖ | | |
| ತಾಲೂಕು | ತಾಲ್ಲೂಕು | | |
| ದಲ್ಲಿನ | ಶಿವಮೊಗ್ಗ | | |
| ಧಾರ್ಮಿಕ | ಆಫ್ರಿಕಾ | | |
| ಗ್ರಂಥ | ಹೊಂದಿದೆ | | |
| ನಿರ್ದೇಶನದ | ಪ್ರಕಾರ | | |
| ವರ್ಗ | ಗೋಕಾಕ | | |
| ಹೊಯ್ಸಳ | ದಿನ | | |
| ಒಬ್ಬ | ಭಾಷೆ | | |

Tamil:

| | | | |
|---|---|---|---|
| நகரில்<br>தஞ்சாவூர்<br>தொழில்கள்<br>மலேசியா<br>சீன<br>இந்தியாவின்<br>இதை<br>ஈரோடு<br>தொடர்பான<br>சிவாலயமாகும் | சிந்து<br>செல்வா<br>விக்கிபீடியா<br>அல்ல<br>கோபி<br>பயனர்கள்<br>பாலாஜி<br>உடைய<br>ஆக்கத்தை<br>அடிப்படையாகக் | பெற்ற<br>சம்பந்தர்<br>அப்பர்<br>சுந்தரர்<br>மூவரதும்<br>பாடல் | வீச்சு<br>என்பது<br>மிக<br>தமிழ்ச்<br>சிங்கப்பூர்<br>பண்பாடு<br>நாட்டார்<br>மாற்றி<br>விடுங்கள்<br>கலைக்களஞ்சியம் |
| நினைக்கிறேன்<br>மிகவும்<br>உலகின்<br>இருக்கும்<br>கூட<br>உலக<br>என்ற<br>தொலைவு<br>தென்னிந்தியாவின்<br>அளவில் | மேலே<br>குறிப்பிட்ட<br>எண்களை<br>ஓடும்<br>சீனாவின்<br>முதன்மையான<br>ஆறுகளில்<br>இருக்கின்றன<br>அறிவித்தது<br>மாற்றங்கள் | ஈடுபாடு<br>மீதும்<br>குறித்த<br>கட்டுரையைத்<br>தொடங்குங்கள்<br>பிற<br>கட்டுரைகளில்<br>தேடிப்பாருங்கள்<br>இலக்கிய<br>கொண்டவர் | நதி<br>உள்ள<br>கோபாலகிருஷ்ணன்<br>உங்களின்<br>த.வி.யில்<br>மொழி<br>பங்கு<br>வகிக்கிறது<br>நன்றி<br>சுந்தர் |
| நீங்கள்<br>அடங்கும்<br>ஸ்ரீநிவாசன்<br>கருத்து<br>உள்ளன<br>ராமசாமி<br>ராவ்<br>என்<br>அமைவது<br>தற்போது | ஒன்று<br>இரண்டு<br>மூன்று<br>நான்கு<br>ஆகிய<br>முக்கிய<br>புகழ்<br>பௌத்த<br>ஆலயம்<br>சிவாஜி<br>கணேசன் | | |

## 2.4 Observation

**(i)      K Means Clustering Error**



**(ii)     K Means Cluster Initialization**

Due to random initialization, one cluster in all the 4 corpora is very big. Few clusters are empty also. But apart from that, words in the remaining clusters are pretty evenly distributed.

**(iii)    K Means Clustering Saturation**

Reduction in error saturates after 4-5 iterations usually.