

# REPORT

# CS B551 Fall 2017, Assignment #3

# Apurva Gupta(guptaapu), Surbhi Paithankar (spaithan), Hasika Mahatta(hmahtta)

# (Based on skeleton code by D. Crandall)

Title: To find POS tags of every word of a given sentence using hidden Markov model.

We have implemented this using three techniques viz.

- a. Simple Naive Bayes Algorithm
- b. Variable Elimination (Forward-Backward algorithm)
- c. Viterbi Algorithm.

a. Simple Naive Bayes Algorithm

Formula used:

$$\text{ARGMAX } P(\text{POSTAG} | \text{WORD } W_i) = P(\text{WORD } W_i | \text{POSTAG}) * P(\text{POSTAG})$$

where  $P(\text{POSTAG})$  is the prior probability for any pos tag

and  $P(\text{WORD } W_i | \text{POSTAG})$  is the emission probability for a word given POSTAG.

For this implementation we use simple bayes net, where every word is independent of other word.

b. Variable Elimination

In Variable elimination we take all other states into account and eliminate them by performing summation over all possible values.

Lets say We need to find  $P(\text{POSTAG } PT_i | \text{WORD } W_i)$

Step 1:

Compute  $\text{tow\_table}(\alpha)$  using forward algorithm from first word upto the desired word. i.e  $PT_1 \dots PT_i$

Step 2:

Compute  $\text{tow\_table}(\beta)$  using backward algorithm from last letter upto the desired letter. i.e  $PT_n \dots PT_i$

Step 3:

Finally we multiply  $\alpha(PT_1 \dots i) * \text{emission}(PT_i | \text{Word } W_i) * \beta(PT_i \dots n)$

c. Viterbi Algorithm

Formula used:

Step 1:

In this we created a viterbi matrix ( $I * J$ ) where  $i$  = number of words in sentence and  $j$  = number of POS tags.

$\text{vit}(i, j)$  stores viterbi values and a backpointer to the previous state.

Formula :

$$\text{vit}(W_0, j) = P(\text{Word } W_0 | \text{POSTAG}) * P(\text{POSTAG})$$

$$\text{vit}(W_i, j) = \max(\text{vit}(W_{i-1}, j) * P(\text{POSTAG} | \text{Prev POSTAG})) * P(\text{Word } W_i | \text{POSTAG})$$

Step 2:

After creating the viterbi matrix, we begin with the last letter and traverse along the path till the first letter using the backpointer maintained.

## PROBABILITY CALCULATIONS:

### 1. Initial Probability:

$P(\text{postag PT}) = \text{Number of occurrences of the postag PT in training file} / \text{total Number of postags}$

Note : If No. of occurrences of POSTAG PT is zero then  $P(\text{Letter L1}) = 10e-14$

### 2. Transition Probability:

$P(\text{POSTAG P2} | \text{POSTAG P1})$  : No. of occurrences where P2 succeeds P1 in training file/ No. of occurrences of P1

Note : If No. of occurrences of POSTAG P1 is zero then  $P(L2 | L1) = 10e-14$

### 3. Emission Probability:

$P(\text{Word } W_i | \text{POSTAG } P_i)$  : No. of occurrences where Word  $W_i$  has POSTAG  $P_i$  in training set/ No. of occurrences of POSTAG  $P_i$ .

## POSTERIOR PROBABILITY:

# Calculate the log of the posterior probability of a given sentence with a given part-of-speech labeling

#The posterior probability is the probability where we need to calculate the probability for each word in the sentence given the corresponding tag for that sentence.

#That is posterior probability =  $P(S1.....Sn | W1.....Wn)$ .

#Applying naive bayes and bayes law, we will get this probability equal to

$P(W1 | S1).....P(Wn | Sn) * P(S1)P(S2/S2).....P(Sn | Sn-1) / P(W1...Wn)$

#(ignoring denominator)

#For example -

#Sentence- The house is big

#label- Det noun pron adj

#In the above bayes net, we can calculate the posterior probability as follows:

$\#P(\text{The/Det}) * P(\text{Det}) * P(\text{house/noun}) * P(\text{noun}) * P(\text{noun/Det}) * P(\text{is/pron}) * P(\text{pron}) * P(\text{pron/noun}) * P(\text{big/adj}) * P(\text{adj}) * P(\text{adj/pron})$

## OBSERVATION:

1. Viterbi & Variable elimination performs much better than simple.

2. Viterbi & VE give almost equal results.

3. We observed that for some really long sentences (like 519th sentence in bc.test set) the issue of underflow takes place. This leads into wrong predication. In order to handle this, we have performed scaling up by multiplying  $10e5$  while computing the tow table. This inturn increased the word accuracy from 94.48% to 95.32% in HMM VE.

## RESULTS:

Scored 2000 sentences with 29442 words.

	Words correct:	Sentences correct:
0. Ground truth:	100.00%	100.00%
1. Simplified:	93.92%	47.45%
2. HMM VE:	95.32%	56.05%
3. HMM MAP:	95.31%	55.30%

