

MIS: 111708049 Assignment Topic - Implement Viterbi Algorithm with backtracking

Files used : Hindi-tagged-[1-12, 14, 15].txt = 14 files

Preprocessing on corpus :

- Remove references of type
- Replace eclipse with space
- Some tags had double underscore, remove the extra underscore
- Separate commas “word\_tag” -> changed to “word\_tag ,\_RD\_PUNC”
- Add RD\_PUNC tag for punctuations where its not present
- Add RD\_SYM tag for symbols where its not present
- Start symbol - “<s>”, End Symbol - “<\s>”
- Adding start and end symbol if not present

Handling sentence ends : Added \_SENT tag for sentence ending with ["?", "!", ".", "l", "|"]

Separating word and tags :

- IF ‘\_’ present : Separate on ‘\_’ to get word and tag associated,
- ELSE :
  - IF word = “<s>”, tag = “START”
  - ELSE IF word = “<\s>”, tag = “END”
  - ELSE : tag = “UN”

Get\_ngrams : returns ngrams using the list of 1-grams(words) provided

Get\_freq\_dict : returns frequency of each ngram in list ngrams

Tag\_Transition\_Matrix(dict) : get all tag bigrams from the list of tags in corpus(sequential),

- IF bigram in corpus :
  - TTP = count(tag bigrams in corpus)/count(tag bigram[0])
- ELSE :
  - TTP = 0.0001

Word Emission probability(dict) : Used add-k smoothing for handling unknown words in test data, |V| = no of distinct words in corpus, get all word-tag pairs,

- IF word-tag pair in corpus :
  - WEP = count(word, tag)+k/count(tag)+k\*V
- ELSE :
  - WEP = k/count(tag)+k\*V

Viterbi Algorithm : return tagged sentence, given test\_sentence, words, tags in corpus(sequential)

Get WEPs, TTPs, SEQSCORE = np.zeros((T, N)), BAKCPTR = np.zeros((T, N))

Initialization Step

For i = 0 to N-1 do

- //N is no of lexical categories and T is no of words
- SEQSCORE (0,i) = P (W 1 |C i ) \* P (C i |”START”)

Iteration Step

For  $t = 1$  to  $T-1$  do

For  $i = 0$  to  $N-1$

$SEQSCORE(t,i) = \text{MAX}_{j=1,N} (SEQSCORE(t-1,j) * P(C_i | C_j)) * P(w_t | C_i)$

$BACKPTR(t,i) = \text{index of } j \text{ that gave max above}$

Sequence Identification Step

$C(T-1) = i$  that maximizes  $SEQSCORE(T-1,i)$

For  $i = T-2$  to  $0$  do

$C(i) = BACKPTR(i+1, C(i+1))$  //Back trace to find the sequence

$tags[C[i]]$  gives tag for word  $test\_tokens[i]$

Preprocess  $test\_sentence$  : Add “<s>”, “</s>” if not present

Main :

Create  $hindi\_tags\_list$ ,  $hindi\_words\_list$  from corpus s.t. Tag of  $hindi\_words\_list[i]$  if given at  $hindi\_tags\_list[i]$  (read sents, preprocess, handle sent ends)

Make TTP, WEP files using respective functions

Get tagged sentences, given unlabelled sentence, and write it to output file.

111708049\_Assign3\_Viterbi\_Input.txt :

आज होगा जबरदस्त मुकाबला क्रिकेट का !

सरकार हमारी जरूरतों को पूरा करे ।

इस नीति के बारे में कुछ भी सही नहीं है ।

आपने क्या उस जगह पर जाने की कोशिश की थी जो मैंने आपको बताया था ?

कई दशक पहले यहां एक बड़ा बरगद का पेड़ हुआ करता था ।

उन्होंने कहा - भारत में आपका स्वागत है !

दिन के अंत तक हम कितनी दूरी तय कर सकते हैं ?

राजा और रानी बहुत घमंडी थे ।

मेरे पास फिल्म को निर्देशित करने के लिए पैसे नहीं हैं ।

नदी के पास एक झोपड़ी है इसलिए हम वहां डेरा डाल सकते हैं !

111708049\_Assign3\_Viterbi\_Output.txt :

आज\_N\_NN होगा\_V\_VM जबरदस्त\_JJ मुकाबला\_N\_NN क्रिकेट\_N\_NN का\_PSP !\_SENT

सरकार\_N\_NN हमारी\_PRP जरूरतों\_N\_NN को\_PSP पूरा\_N\_NN करे\_V\_VM ।\_SENT

इस\_DM\_DMD नीति\_N\_NN के\_PSP बारे\_N\_NN में\_PSP कुछ\_QT\_QTF भी\_RP\_RPD

सही\_N\_NN नहीं\_RP\_NEG है\_V\_VAUX ।\_SENT

आपने\_PRP क्या\_PRQ उस\_PRP जगह\_N\_NN पर\_PSP जाने\_V\_VM की\_PSP कोशिश\_N\_NN

की\_PSP थी\_N\_NN जो\_PRP मैंने\_PRP आपको\_PRP बताया\_V\_VM था\_V\_VAUX ?\_SENT

कई\_QT\_QTF दशक\_N\_NN पहले\_IN यहां\_N\_NN एक\_QT\_QTC बड़ा\_JJ बरगद\_N\_NN

का\_PSP पेड़\_N\_NN हुआ\_V\_VM करता\_V\_VM था\_V\_VAUX ।\_SENT

उन्होंने\_N\_NN कहा\_V\_VM -\_RD\_SYM भारत\_N\_NN में\_PSP आपका\_PRP स्वागत\_N\_NN

है\_V\_VM !\_SENT

दिन\_N\_NN के\_PSP अंत\_N\_NN तक\_PSP हम\_PRP कितनी\_QTF दूरी\_N\_NN तय\_N\_NN  
कर\_V\_VM सकते\_V\_VAUX हैं\_V\_VAUX ?\_SENT  
राजा\_N\_NN और\_CCD रानी\_N\_NN बहुत\_PSP घमंडी\_V\_VM थे\_V\_VAUX |\_SENT  
मेरे\_PR\_PRP पास\_N\_NST फिल्म\_N\_NN को\_PSP निर्देशित\_N\_NN करने\_V\_VM के\_PSP  
लिए\_PSP पैसे\_N\_NN नहीं\_RP\_NEG हैं\_V\_VAUX |\_SENT  
नदी\_N\_NN के\_PSP पास\_N\_NST एक\_QT\_QTC झोपड़ी\_N\_NN है\_V\_VB इसलिए\_CC  
हम\_PRP वहां\_N\_NN डेरा\_N\_NN डाल\_V\_VM सकते\_V\_VAUX हैं\_V\_VAUX !\_SENT