

# POS TAGGING

Using Hidden Markov Model

By

Ravi Yadav

Biplab Ch Das

# Generating the Lexical Probabilities

We know that  $p(W | T) = P(W, T) / P(T)$

And since  $p(W, T) = N(W, T) / N(T)$

We Used a Hash Map to get all the counts of tagged words in the training corpus. That is we first find  $N(W, T)$ . And then divide it by  $N(T)$ .



# Generating The transition Probabilities

We know that  $p(T_n | T_{n-1}) = P(T_n, T_{n-1}) / P(T_{n-1})$

And since  $p(W, T) = N(T_{n-1}, T_n) / N(T_{n-1})$

Used Hash Map again. Got the  $N(t(n-1), t(n))$  values and then we divided it by  $N(t(n-1))$ .

# Separation of folds

The folds were separated in the ratio of 4:1  
4 for training and 1 for the testing.

5 folds were created and the models were tested  
and trained on each folds.

# Smoothing of the data

For the non existent data we tried to smoothen the data by adding a minor value(0.0000001) to the count.

A single zero in product can make the whole term zero.

# The viterbi algorithm

The viterbi algorithm was implemented to get the best possible state sequence using the transition matrix and the emission matrix. The best possible state sequence was considered as the best tag possible for the sentence

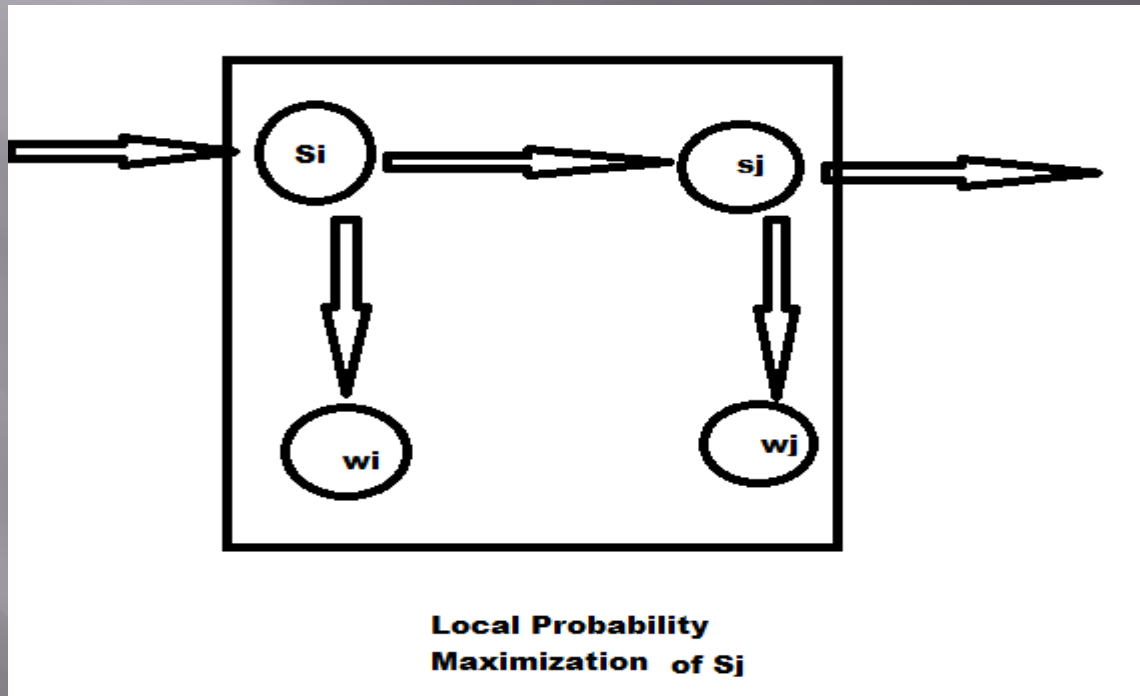
# The previous algorithm implemented

Till now the pos tagger have implemented the viterbi algorithm ,at the same time it we kept the greedy local lookup that we presented before.

$$T(n)^* = \operatorname{argmax}_{(T_j)} (p((T_j | T_i)^* p(W_j | t_j));$$

Here 'i' is prev state and 'j' varies.

# Belief network of assumption



We take joint probability

$$P(S_i, S_j, w_i, w_j) = P(s_j | s_i) * P(s_i) * P(w_j | s_j) * P(w_i | s_i)$$

And set  $P(s_i) = 1$  and  $P(w_i | s_i) = 1$  (Considering the prev state is constant, we don't need to consider  $P(s_i)$  for local maximisation)

# Confusion Matix

The confusion matrix have been created. The actual tags were shown in the right hand side column and the tags by the Viterbi algorithm were enumerated at the bottom line.

[file:///C:/Users/Biplab/Desktop/confusion\\_matrix.html](file:///C:/Users/Biplab/Desktop/confusion_matrix.html)



# Confusion Reason

Most of the confusion was between NP0 and NN1  
It was because of the unknown word prediction  
handling done in the code.

# The F-score

The proper output for the tag wise precision-recall and the F-score has been calculated .

The F-score was calculated putting the beta value to 1.

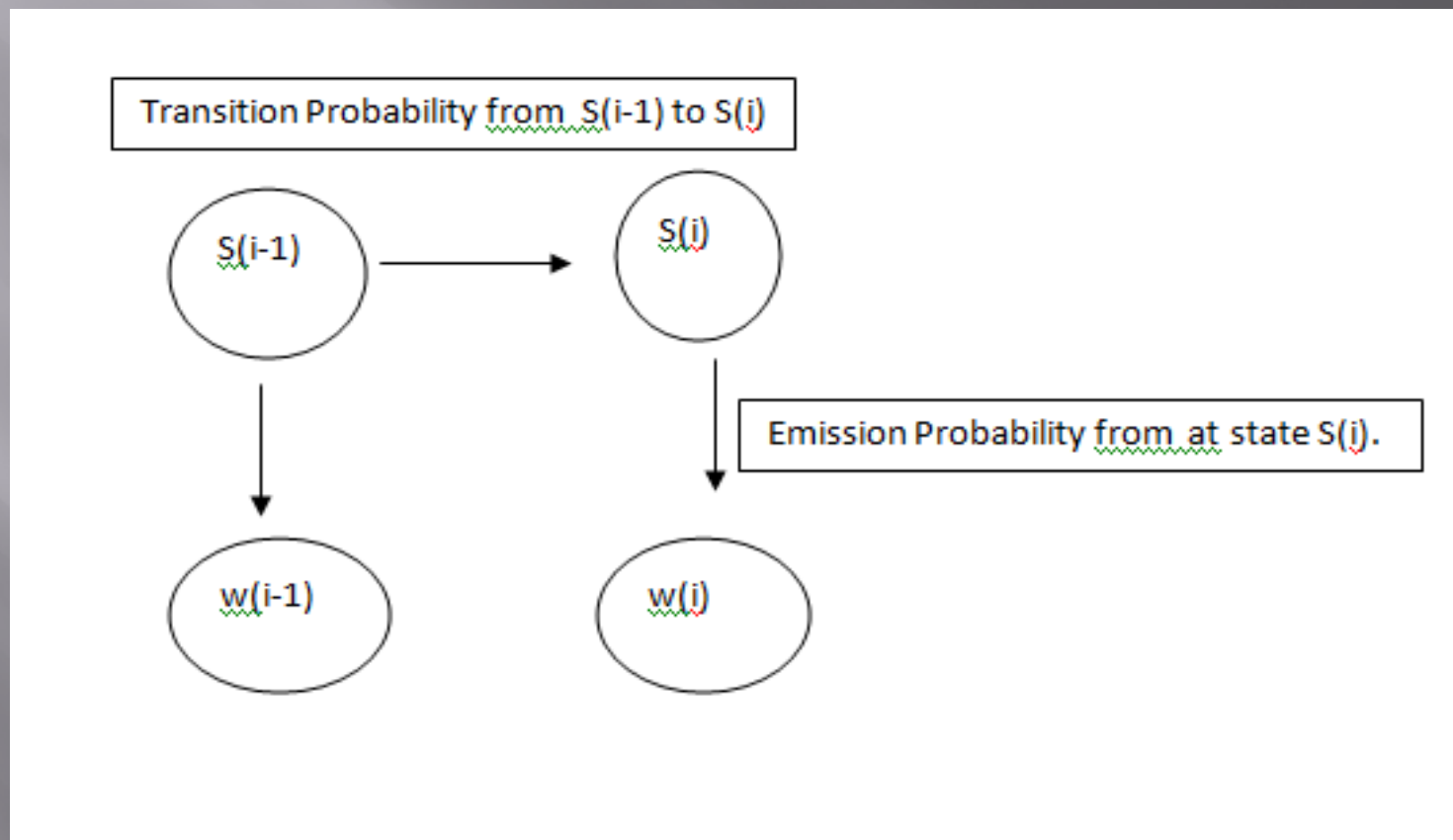
<file:///C:/Users/Biplab/Desktop/proper.html>

# Trigram Implentation

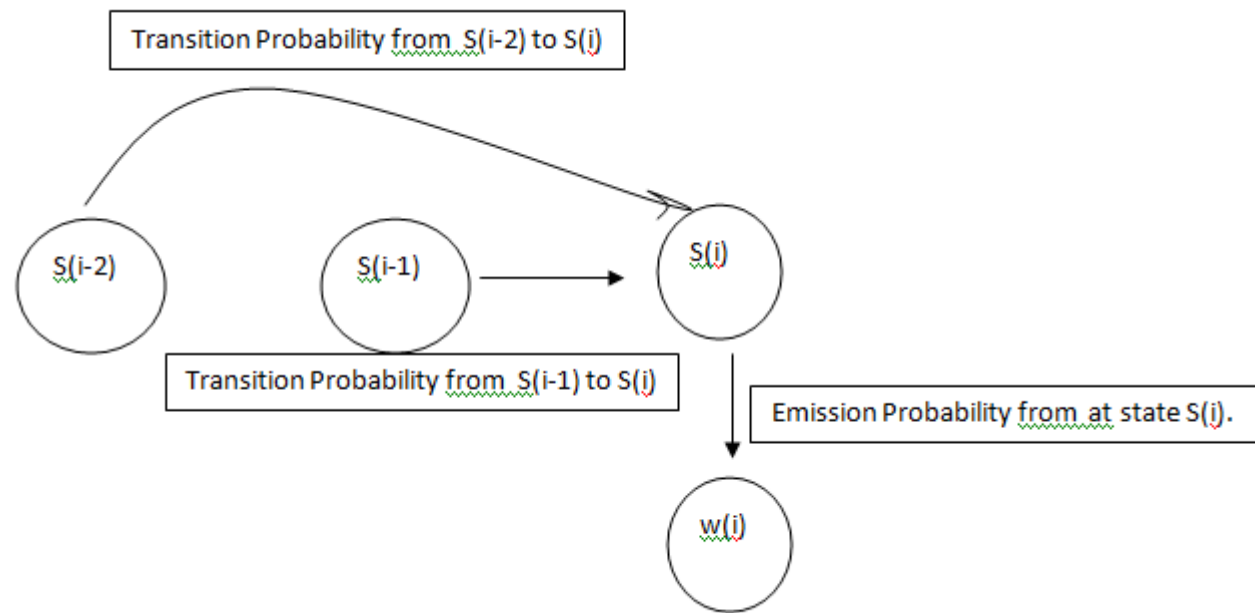
The trigram model was implemented for the Pos tagger.

In this case the current tag not only depends on the previous state but also on the previous to previous state.

# Trigram the extension of bigram model:



# The graphical model for trigram model:



# A star algorithm

Heuristic used:

$$-h_{(Ti)} = \log(P(W_i | T_i)) + \sum_{(j=i+1 \text{ to } n)} \log(\max(P(W_j | T_j))) + \sum_j \log(\max(T_i | T_j))$$

Where,

$P(W_i | T_i)$  is the lexical prob at a node.

$\max(T_i | T_j)$  is the maximum of all transitional probabilities

$\max(T_i | T_j)$  is the max lexical prob at a node

# Accuracy of the tagger

Fold No	Bigram	Skip gram ( $0.75 \times \text{bigram} + 0.25 \times \text{skipgram}$ )	( $0.23 \times \text{bigram} + 0.11 \times \text{skipgram} + 0.66 \times \text{trigram}$ )	A* algorithm
Fold 1	92.69	92.71	92.82	93.20
Fold 2	93.76	94.08	94.28	94.75
Fold 3	93.61	93.69	93.73	94.66
Fold 4	94.00	94.06	94.17	94.72
Fold 5	93.22	93.38	93.64	93.86
Overall	93.45	93.58	93.73	94.23



# Discriminative

Model tag directly from the data.

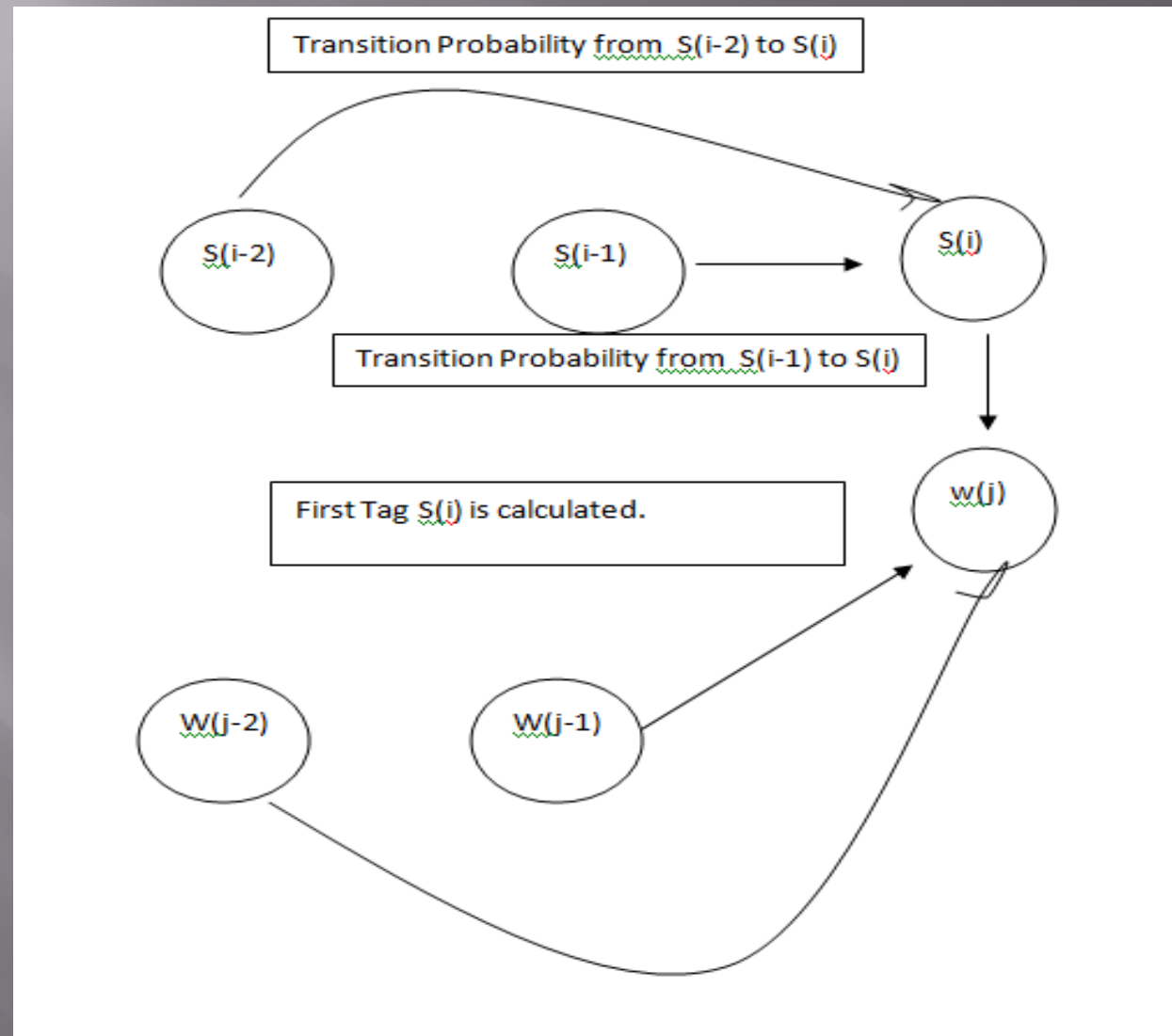
Expression :

$\text{Argmax}(j) P(T(j) | T(j-1), \text{Current\_word})$

Accuracy achieved 72.98 %

Possible problem : Unknown word handling was not done .

# Word Prediction



# Model for word prediction:

$$\operatorname{argmax}_{Nw} \sum_{Tw} I(Tw) * \operatorname{sim}(\operatorname{Syn}(Nw), \operatorname{Syn}(Tw))$$

Where  $\operatorname{argmax}$  is over  $Nw$

Here  $nw = \operatorname{argmax}(w_i, w_j)$  maximizing over  $w_j$

$\operatorname{Sim}$  is a similarity function

$\operatorname{Syn}(W)$  refers to the synset of the words

# Accuracy of the word prediction model

Measured using trigram of content words and checked with the top 5 suggestions.

It came out to be 43.30% without knowing pos tags and with pos tags it increased to 46.00%

Most of the confusion occurred when the next words involved names.