

REPORT P2-A

In the assignment we are supposed to find out the Part of Speech tags using Hidden Markov Model. The brown corpus is used for training and testing the Model.

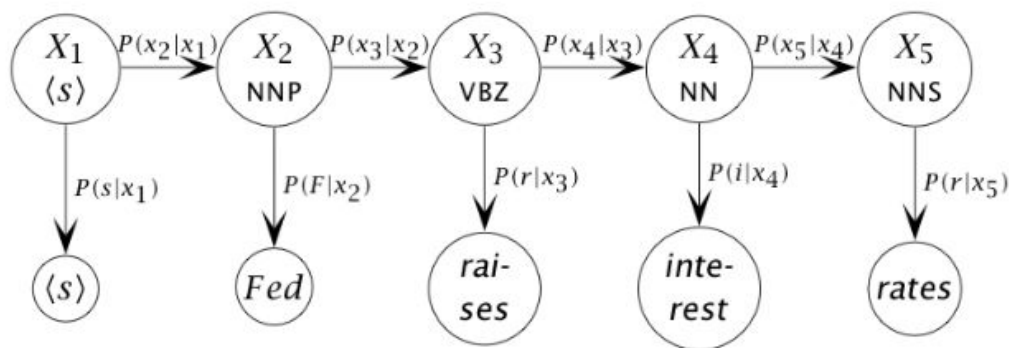
The training data is in the form of $[(x_1, y_1), (x_2, y_2), \dots]$ where x_i is the word in the sentence and y_i is the corresponding tag to the word.

Hidden Markov Models (HMMs) model : A system of discrete temporal unobserved (hidden) variables and discrete temporal observed variables. The observed and unobserved variables are related through emission probabilities.

Viterbi algorithm : Finding most likely sequence in HMM

The Viterbi algorithm is a dynamic programming algorithm that efficiently computes the the most likely states of the latent variables of a Hidden Markov Model (HMM), given an observed sequence of emissions.

In tagging problem, each x_i would be a sequence of words $x_1 x_2 x_3 \dots x_n$ and each y_i would be a sequence of tags $y_1 y_2 y_3 \dots y_n$



Implementation Details :

1. Pre-processing

2. Data Structure

- A dictionary is maintained to store the trigrams, bigrams, state transition and states.

(Here tags are the states and words are transition).

- Smoothing is used so as to handle the unknown words.