

Name : Sarvagya Singh

SAPID : 60009200030

BATCH : K1

EXPERIMENT

NO 2

AIM: To make understand and learn the technique and use of N-Gram Model

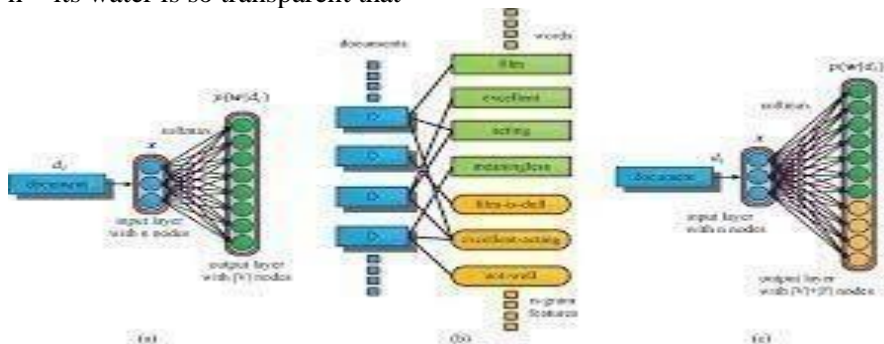
THEORY:

Statistical language models, in its essence, are the type of models that assign probabilities to the sequences of words. In this article, we'll understand the simplest model that assigns probabilities to sentences and sequences of words, the n-gram

You can think of an N-gram as the sequence of N words, by that notion, a 2-gram (or bigram) is a two-word sequence of words like "please turn", "turn your", or "your homework", and a 3-gram (or trigram) is a three-word sequence of words like "please turn your", or "turn your homework"

Let's start with equation $P(w|h)$, the probability of word w , given some history, h . For example, Here, w = The word given

h = its water is so transparent that



And, one way to estimate the above probability function is through the relative frequency count approach, where you would take a substantially large corpus, count the number of times you see its water is so transparent that, and then count the number of times it is followed by the. In other words, you are answering the question: Out of the times you saw the history h , how many times did the word w follow it. Now, you can imagine it is not feasible to perform this over an entire corpus; especially it is of a significant size.

This shortcoming and ways to decompose the probability function using the chain rule serves as the base intuition of the N-gram model. Here, you, instead of computing probability using the entire corpus, would approximate it by just a few historical words

The Bigram Model

As the name suggests, the bigram model approximates the probability of a word given all the previous words by using only the conditional probability of one preceding word. In other words, you approximate it with the probability: $P(\text{the} | \text{that})$

And so, when you use a bigram model to predict the conditional probability of the next word, you are thus making the following approximation:

This assumption that the probability of a word depends only on the previous word is also known as **Markov assumption**.

Markov models are the class of probabilistic models that assume that we can predict the probability of some future unit without looking too far in the past.

Name : Sarvagya Singh

SAPID : 60009200030

BATCH : K1

You can further generalize the bigram model to the trigram model which looks two words into the past and can thus be further generalized to the N-gram model.

Data Set to Be Used

Lab Experiments to be performed in this Session: Step-by-step implementation of ngram language model

1. Basic pre-processing
2. Code to generate N-grams
3. Creating unigrams
4. Creating bigrams
5. Creating trigrams
6. Finding Frequency Distribution
7. Finding Probabilities for Bigram
8. Finding Next word for the given word using MLE (Maximum Likelihood Estimate)

