 <b>Marwadi University</b>	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence (01CT0616)</b>	<b>Aim:</b> To study N-gram model	
<b>Experiment No: 7</b>	<b>Date:</b>	<b>Enrolment No:</b>

**Aim:** To study N-gram model

**IDE:** Google Colab

### **Theory:**

Given a sequence of N-1 words, an N-gram model predicts the most probable word that might follow this sequence. It's a probabilistic model that's trained on a corpus of text. Such a model is useful in many NLP applications including speech recognition, machine translation and predictive text input.

An N-gram model is built by counting how often word sequences occur in corpus text and then estimating the probabilities. Since a simple N-gram model has limitations, improvements are often made via smoothing, interpolation and backoff. An N-gram model is one type of a Language Model (LM), which is about finding the probability distribution over word sequences.

Consider two sentences: "There was heavy rain" vs. "There was heavy flood". From experience, we know that the former sentence sounds better. An N-gram model will tell us that "heavy rain" occurs much more often than "heavy flood" in the training corpus. Thus, the first sentence is more probable and will be selected by the model.

A model that simply relies on how often a word occurs without looking at previous words is called unigram. If a model considers only the previous word to predict the current word, then it's called bigram. If two previous words are considered, then it's a trigram model.

An n-gram model for the above example would calculate the following probability:


$$P(\text{'There was heavy rain'}) = P(\text{'There', 'was', 'heavy', 'rain'}) = P(\text{'There'})P(\text{'was' | 'There'})P(\text{'heavy' | 'There was'})P(\text{'rain' | 'There was heavy'})$$

Since it's impractical to calculate these conditional probabilities, using Markov assumption, we approximate this to a bigram model:  $P(\text{'There was heavy rain'}) \sim P(\text{'There'})P(\text{'was' | 'There'})P(\text{'heavy' | 'was'})P(\text{'rain' | 'heavy'})$  In speech recognition, input may be noisy and this can lead to wrong speech-to-text conversions.

N-gram models can correct this based on their knowledge of the probabilities. Likewise, N-gram models are used in machine translation to produce more natural sentences in the target language.

When correcting for spelling errors, sometimes dictionary lookups will not help. For example, in the phrase "in about fifteen minuets" the word 'minuets' is a valid dictionary word but it's incorrect in this context. N-gram models can correct such errors.

N-gram models are usually at word level. It's also been used at character level to do stemming, that is, separate the root word from the suffix. By looking at N-gram statistics, we could also classify languages or differentiate

 <b>Marwadi</b> University	<b>Marwadi University</b> <b>Faculty of Technology</b> <b>Department of Information and Communication Technology</b>	
<b>Subject: Artificial Intelligence (01CT0616)</b>	<b>Aim:</b> To study N-gram model	
<b>Experiment No: 7</b>	<b>Date:</b>	<b>Enrolment No:</b>

between US and UK spellings. For example, 'sz' is common in Czech; 'gb' and 'kp' are common in Igbo. In general, many NLP applications benefit from N-gram models including part-of-speech tagging, natural language generation, word similarity, sentiment extraction and predictive text input.

To preprocess your text simply means to bring your text into a form that is predictable and analyzable for your task. A task here is a combination of approach and domain. Machine Learning needs data in the numeric form. We basically used encoding technique (Bag Of Word, Bi-gram, n-gram, TF-IDF, Word2Vec) to encode text into numeric vector. But before encoding we first need to clean the text data and this process to prepare (or clean) text data before encoding is called text preprocessing, this is the very first step to solve the NLP problems.

**Program (Code):**

- To be attached with
1. N-Gram Model Code for generating dynamic N-Grams

**Results:**

To be attached with

**Observation:**

**Post Lab Exercise:**

Apply the probability distribution concept to predict the next word by taking sample document and seed phrases. Attach the code and screenshot. Write your observation here.

---



---



---



---