

Tribhuvan University  
Institute of Engineering  
Pulchowk Campus

---

A Project Report on:  
**Nepali Language Processing**

---

**Submitted By:**

Amrit Baral(076bct006)  
Nabin Da Shrestha(076bct037)  
Nirajan Bekoju(076bct039)  
Nishant Luitel (076bct041)

**Submitted To:**

Department of Electronics and Computer Engineering

**Submission Date:**

8<sup>th</sup> November, 2022

# Contents

<b>1</b>	<b>Probabilistic Language Model</b>	<b>5</b>
1.1	Introduction . . . . .	5
1.2	N-gram . . . . .	6
1.3	Sequence Notation . . . . .	6
1.4	N-gram Probability . . . . .	6
1.4.1	Uni-gram Probability . . . . .	6
1.4.2	Bi-gram Probability . . . . .	7
1.4.3	Tri-gram Probability . . . . .	7
1.5	Probability of a sequence . . . . .	8
1.6	Approximation of Sequence Probability . . . . .	8
1.7	Starting and Ending Sentences . . . . .	9
1.7.1	Start of sentence symbol $\langle s \rangle$ . . . . .	9
1.7.2	End of sentence symbol $\langle /s \rangle$ . . . . .	9
1.8	Count Matrix . . . . .	10
1.9	Probability Matrix . . . . .	10
1.10	Generative Language Model . . . . .	11
1.11	Train, Validation and Test Split . . . . .	12
1.12	Language Model Evaluation . . . . .	12
1.12.1	Extrinsic evaluation . . . . .	12
1.12.2	Intrinsic evaluation . . . . .	12
1.12.3	Perplexity . . . . .	13
1.12.4	Log Perplexity . . . . .	13
1.13	Vocabulary . . . . .	14
1.14	Missing N-gram in training corpus . . . . .	14
1.14.1	Laplacian Smoothing . . . . .	15
1.14.2	K-smoothing . . . . .	15

# List of Figures

1.1	Language Model for Auto-complete . . . . .	5
-----	--	---

# List of Tables

1.1	Count Matrix for bi-gram model . . . . .	10
1.2	Count Matrix for bi-gram model with row-sum . . . . .	10
1.3	Probability Matrix for bi-gram model . . . . .	11
1.4	Train, Validation and Test Split . . . . .	12

# Chapter 1

## Probabilistic Language Model

### 1.1 Introduction

Language model is used to estimate the probability of the word sequences and to estimate the probability of a word following the sequence of words. The concept can be used to auto-complete a sentence with most likely suggestions as shown below:

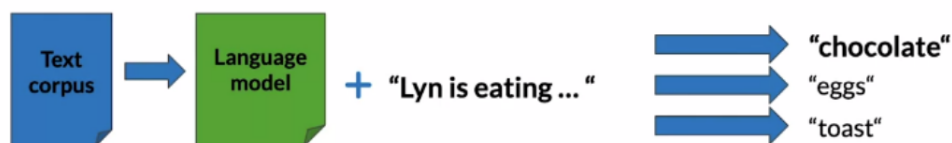


Figure 1.1: Language Model for Auto-complete

Some other applications of the language model are:

1. Speech Recognition
2. Spelling Correction
3. Augmentative Communication

## 1.2 N-gram

A N-gram is a sequence of N-words. For example:

**Corpus :** I am happy because I am learning.

**Uni-grams :** {I, am, happy, because, I, am, learning }

**Bi-grams :** {I am, am happy, happy because, ... }

**Tri-grams :** {I am happy, am happy because, ... }

## 1.3 Sequence Notation

Let's define a standard sequence notation to represent the sequence of words in our corpus. Let our corpus be as shown below where  $n^{th}$  is represented as  $W_n$ . Let the number of words in the corpus be  $m = 500$

**Corpus:** This is great. ... teacher drinks tea.

Hence,  $W_1 = \text{This}$ ,  $W_2 = \text{is}$ ,  $W_3 = \text{great}$  and so on...

Let us represent the sequence of words as  $W_s^e$  where  $s$  represent the starting index and  $e$  represent the ending index of the words in the corpus. So,

$W_1^m = W_1 W_2 W_3 \dots W_{m-1} W_m$

Similarly,  $W_1^3 = W_1 W_2 W_3 = \text{This is great}$

## 1.4 N-gram Probability

N-gram probability is the probability of the occurrence of the sequence of N-words in a corpus. From uni-gram, bi-gram and tri-gram probability, we can deduce for N-gram probability denoted by  $P(W_N|W_1^{N-1})$  as

$$P(W_N|W_1^{N-1}) = \frac{C(W_1^{N-1} W_N)}{C(W_1^{N-1})} = \frac{C(W_1^N)}{C(W_1^{N-1})} \quad (1.1)$$

### 1.4.1 Uni-gram Probability

Uni-gram Probability is the probability of occurrence of a word in a given corpus.

**Corpus :** I am happy because I am learning.

**Size of corpus :**  $m = 7$

$$P(I) = \frac{2}{7} \quad (1.2)$$

$$P(happy) = \frac{1}{7} \quad (1.3)$$

Hence, we can see that Probability of uni-gram is given by:

$$P(w) = \frac{C(w)}{m} \quad (1.4)$$

where  $C(w)$  = frequency of the word in the corpus

### 1.4.2 Bi-gram Probability

Bi-gram Probability is the probability of occurrence of sequence of two words in a corpus.

**Corpus :** I am happy because I am learning.

**Size of corpus :**  $m = 7$

$$P(am|I) = \frac{C(I\ am)}{C(I)} = \frac{2}{2} = 1 \quad (1.5)$$

$$P(happy|I) = \frac{C(I\ happy)}{C(I)} = \frac{0}{7} = 0 \quad (1.6)$$

$$P(learning|am) = \frac{C(am\ learning)}{C(am)} = 1/2 \quad (1.7)$$

Hence, we can see that Probability of bi-gram is given by:

$$P(y|x) = \frac{C(x\ y)}{C(x)} \quad (1.8)$$

where  $C(w)$  = frequency of the word in the corpus

### 1.4.3 Tri-gram Probability

Tri-gram Probability is the probability of the occurrence of the sequence of three words in a corpus.

**Corpus :** I am happy because I am learning.

**Size of corpus :**  $m = 7$

$$P(happy|I\ am) = \frac{C(I\ am\ happy)}{C(I\ am)} = \frac{1}{2} \quad (1.9)$$

Hence, we can see that Probability of bi-gram is given by:

$$P(W_3|W_1^2) = \frac{C(W_1^2 W_3)}{C(W_1^2)} = \frac{C(W_1^3)}{C(W_1^2)} \quad (1.10)$$

where  $C(w)$  = frequency of the word in the corpus

## 1.5 Probability of a sequence

From the conditional probability and chain rule, we have

$$P(B|A) = \frac{P(A, B)}{P(A)} \quad (1.11)$$

From (1.11), we get

$$P(A, B) = P(A).P(B|A) \quad (1.12)$$

Using (1.12), we can deduce

$$P(A, B, C, D) = P(A).P(B|A).P(C|A, B).P(D|A, B, C) \quad (1.13)$$

Now, the probability of the sequence of words: "the teacher drinks tea" is given by:

$$\begin{aligned} P(\text{the teacher drinks tea}) &= P(\text{the}).P(\text{teacher}|\text{the}) \\ &P(\text{drinks}|\text{the teacher}).P(\text{tea}|\text{the teacher drinks}) \end{aligned} \quad (1.14)$$

## 1.6 Approximation of Sequence Probability

When using n-gram language model, the exact sentence for the required n-gram might not be present exactly in corpus due to which the frequency of the words sequence becomes zero and hence the probability. For example:

**Input :** the teacher drinks tea

$P(\text{the teacher drinks tea}) = P(\text{the}) \cdot P(\text{teacher} | \text{the}) P(\text{drinks} | \text{the teacher})$   
 $\cdot P(\text{tea} | \text{the teacher drinks})$

Here,

$$P(\text{tea}|\text{the teacher drinks}) = \frac{C(\text{the teacher drinks tea})}{C(\text{the teacher drinks})} \quad (1.15)$$



In the above equation (1.15), the frequency of "the teacher drinks tea" and "the teacher drinks" are both likely 0.

Hence, in order to calculate the required probability of the sequence of the words, as per Markov assumption, only last N words matter. Suppose, only last one words matter then, we have

$$\begin{aligned} & P(\text{the teacher drinks tea}) \\ &= P(\text{the}) \cdot P(\text{teacher} \mid \text{the}) P(\text{drinks} \mid \text{the teacher}) \cdot P(\text{tea} \mid \text{the teacher drinks}) \\ &\approx P(\text{the}) \cdot P(\text{teacher} \mid \text{the}) P(\text{drinks} \mid \text{teacher}) \cdot P(\text{tea} \mid \text{drinks}) \end{aligned}$$

Hence, we can model entire sentence with n-gram following Markov Assumption as follow

$$P(W_n | W_1^{n-1}) \approx P(W_n | W_{n-N+1}^{n-1}) \quad (1.16)$$

## 1.7 Starting and Ending Sentences

### 1.7.1 Start of sentence symbol $\langle s \rangle$

For an input sentence, "the teacher drinks tea", we have

$$P(\text{the teacher drinks tea}) \approx P(\text{the}) \cdot P(\text{teacher} \mid \text{the}) P(\text{drinks} \mid \text{teacher}) \cdot P(\text{tea} \mid \text{drinks})$$

After, we add a start token  $\langle s \rangle$ , we have

**Input :**  $\langle s \rangle$  the teacher drinks tea

$$P(\text{the teacher drinks tea}) \approx P(\text{the} \mid \langle s \rangle) \cdot P(\text{teacher} \mid \text{the}) P(\text{drinks} \mid \text{teacher}) \cdot P(\text{tea} \mid \text{drinks})$$

In this way, we can add start token to determine the possible starting words and calculate the probability of sequence considering first word in N-gram language model. For Tri-gram and N-gram, we add (N - 1) start token  $\langle s \rangle$  in the beginning of the sequence.

### 1.7.2 End of sentence symbol $\langle /s \rangle$

For an input sentence, "the teacher drinks tea", we have

$$P(\text{the teacher drinks tea}) \approx P(\text{the}) \cdot P(\text{teacher} \mid \text{the}) P(\text{drinks} \mid \text{teacher}) \cdot P(\text{tea} \mid \text{drinks})$$

After, we add a start token  $\langle s \rangle$  and an end token  $\langle /s \rangle$ , we have

**Input :**  $\langle s \rangle$  the teacher drinks tea  $\langle /s \rangle$

$P(\text{the teacher drinks tea}) \approx P(\text{the} | \langle s \rangle) \cdot P(\text{teacher} | \text{the}) P(\text{drinks} | \text{teacher})$   
 $\cdot P(\text{tea} | \text{drinks}) \cdot P(\langle /s \rangle | \text{tea})$

For bi-gram, tri-gram or any N-gram model, we add end token just once at the end of sentence unlike start token

## 1.8 Count Matrix

Count matrix is the matrix containing the frequency of occurrence of N-gram. Rows in count matrix represent the unique corpus of (N - 1) gram and columns represent the unique corpus word.

**Corpus :**  $\langle s \rangle$  I study I learn  $\langle /s \rangle$

Bi-gram Count Matrix is as follow:

	$\langle s \rangle$	$\langle /s \rangle$	I	study	learn
$\langle s \rangle$	0	0	1	0	0
$\langle /s \rangle$	0	0	0	0	0
I	0	0	0	1	1
study	0	0	1	0	0
learn	0	1	0	0	0

Table 1.1: Count Matrix for bi-gram model

## 1.9 Probability Matrix

Probability matrix is the matrix containing the probability of the occurrence of N-gram words sequence. It can be created by dividing each elements of the count matrix by total sum of all row elements. For example:

**Corpus :**  $\langle s \rangle$  I study I learn  $\langle /s \rangle$

	$\langle s \rangle$	$\langle /s \rangle$	I	study	learn	sum
$\langle s \rangle$	0	0	1	0	0	1
$\langle /s \rangle$	0	0	0	0	0	0
I	0	0	0	1	1	2
study	0	0	1	0	0	1
learn	0	1	0	0	0	1

Table 1.2: Count Matrix for bi-gram model with row-sum

	$\langle s \rangle$	$\langle /s \rangle$	I	study	learn
$\langle s \rangle$	0	0	1	0	0
$\langle /s \rangle$	0	0	0	0	0
I	0	0	0	0.5	0.5
study	0	0	1	0	0
learn	0	1	0	0	0

Table 1.3: Probability Matrix for bi-gram model

From the above probability matrix, we can calculate the probability of sentence and also predict the next word based on higher probability.

For example :

**Input :**  $\langle s \rangle$  I learn  $\langle /s \rangle$

$$\begin{aligned}
P(\text{Input}) &= P(\langle s \rangle \text{ I learn } \langle /s \rangle) \\
&= P(I \mid \langle s \rangle) \cdot P(\text{learn} \mid I) \cdot P(\langle /s \rangle \mid \text{learn}) \\
&= 1 \times 0.5 \times 1 \\
&= 0.5
\end{aligned}$$

## 1.10 Generative Language Model

Generative language model is used to generate a most likely sentences.

Algorithm for the generative language model is:

Step 1 : Choose sentence start

Step 2 : Choose next bi-gram starting with the previous word

Step 3 : Continue until  $\langle /s \rangle$  is picked

For example:

Assume our corpus be:

$\langle s \rangle$  Lyn drinks chocolate  $\langle /s \rangle$

$\langle s \rangle$  John drinks tea  $\langle /s \rangle$

$\langle s \rangle$  Lyn eats chocolate  $\langle /s \rangle$

Then, the generative language model can generate a sentence as follow:

Step 1: ( $\langle s \rangle$ , Lyn) or ( $\langle s \rangle$ , John)

Step 2: (Lyn, eats) or (Lyn drinks)

Step 3: (drinks, tea) or (drinks, chocolate)

Step 4: (tea,  $\langle /s \rangle$ )

## 1.11 Train, Validation and Test Split

The recommended split of corpus into train, validation and test data is as follow

Type of corpus	Train	Validation	Test
Small	80%	10%	10%
Large	98%	1%	1%

Table 1.4: Train, Validation and Test Split

## 1.12 Language Model Evaluation

We can use two different approaches to evaluate and compare language models:

### 1.12.1 Extrinsic evaluation

This involves evaluating the models by employing them in an actual task (such as machine translation) and looking at their final loss/accuracy. This is the best option as it's the only way to tangibly see how different models affect the task we're interested in. However, it can be computationally expensive and slow as it requires training a full system.

### 1.12.2 Intrinsic evaluation

This involves finding some metric to evaluate the language model itself, not taking into account the specific tasks it's going to be used for. While intrinsic evaluation is not as “good” as extrinsic evaluation as a final metric, it's a useful way of quickly comparing models. Perplexity is an intrinsic evaluation method.

### 1.12.3 Perplexity

Perplexity is an evaluation metric for language models. For better language model, perplexity value should be smaller. Perplexity can be calculated as follow:

$$PP(W) = P(s_1, s_2, s_3, s_4, \dots, s_m)^{\frac{-1}{m}} \quad (1.17)$$

where,

$\mathbf{W}$  = test set containing  $m$  sentences  $s$

$S_i = i^{th}$  sentence in the test set, each ending with  $\langle /s \rangle$

$\mathbf{m}$  = number of all words in entire test set  $\mathbf{W}$  including  $\langle /s \rangle$  but not including  $\langle s \rangle$

Similarly, Perplexity for the bi-gram model can be calculated as:

$$PP(W) = \sqrt[m]{\prod_{i=1}^m \prod_{j=1}^{|S_i|} \frac{1}{P(w_j^{(i)} | w_{j-1}^{(i)})}} \quad (1.18)$$

where,  $w_j^{(i)} = j^{th}$  word in  $i^{th}$  sentence

If we concatenate all the sentences in  $\mathbf{W}$ , we can get perplexity as

$$PP(W) = \sqrt[m]{\prod_{i=1}^m \frac{1}{P(w_i | w_{i-1})}} \quad (1.19)$$

where,  $w_i = i^{th}$  word in test set

### 1.12.4 Log Perplexity

Log perplexity is calculated by taking logarithm on both sides of equation (1.19) and we get

$$\log PP(W) = \frac{-1}{m} \sum_{i=1}^m \log_2(P(w_i | w_{i-1})) \quad (1.20)$$

## 1.13 Vocabulary

Vocabulary is a set of words from the training corpus. We can create vocabulary from training corpus using following criteria.

1. Min word frequency  $f$
2. Max  $|V|$ , include words by frequency
3. Using  $\langle UNK \rangle$  for unknown token

For example: Suppose our corpus is as follow

$\langle s \rangle$  Lyn drinks chocolate  $\langle /s \rangle$

$\langle s \rangle$  John drinks tea  $\langle /s \rangle$

$\langle s \rangle$  Lyn eats chocolate  $\langle /s \rangle$

let min frequency to be added in vocabulary be  $f = 2$ . So, now corpus will be as follow:

$\langle s \rangle$  Lyn drinks chocolate  $\langle /s \rangle$

$\langle s \rangle$   $\langle UNK \rangle$  drinks  $\langle UNK \rangle$   $\langle /s \rangle$

$\langle s \rangle$  Lyn  $\langle UNK \rangle$  chocolate  $\langle /s \rangle$

Hence, vocabulary =  $\{Lyn, drinks, chocolate\}$

Suppose,

**Input Query :**  $\langle s \rangle$  Adam drinks chocolate  $\langle /s \rangle$

So, the input query after processing becomes

**Ouput :**  $\langle s \rangle$   $\langle UNK \rangle$  drinks chocolate  $\langle /s \rangle$

In this way, vocabulary can be created and unknown words can be handled.

## 1.14 Missing N-gram in training corpus

N-grams made of known words still might be missing in the training corpus. N-gram probability is given as:

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}, w_n)}{C(w_{n-N+1}^{n-1})} \quad (1.21)$$

In equation (1.21), both  $C(w_{n-N+1}^{n-1}, w_n)$  and  $C(w_{n-N+1}^{n-1})$  can be 0 in the probability matrix which causes the probability to be indefinite. It can be solved by using smoothing.

### 1.14.1 Laplacian Smoothing

In Laplacian smoothing, probability is calculated as follow:

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}, w_n) + 1}{C(w_{n-1}) + V} \quad (1.22)$$

### 1.14.2 K-smoothing

In K-smoothing, probability is calculated as follow

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}, w_n) + K}{C(w_{n-1}) + k * V} \quad (1.23)$$