

Probability and Statistics

MA – 202, Project 1

# Statistical Language Modelling Using N-grams

MENTOR

Prof. Mayank Singh



# I Project Brief

Our aim, using Natural Language Processing, is to predict the most probable next word and analyse the correctness of the input sentence. Further, we would also create a User interface to do the same to make it more interactive and user friendly.





We divided the task into the following main steps and then integrated it to get the final result:

- Corpus Preparation
- Laplace Smoothing
- Next word prediction
- Perplexity of sentence
- Interface



## Reliable Dataset:

- The model's accuracy depends on the reliability of the dataset used to train the model.
- For the given model, we have sought and utilized the best dataset available in our view

## Data preprocessing:

- Firstly, we downloaded a suitable parser to convert the acquired dataset into machine language.
- Then, we installed suitable libraries to enhance preprocessing of data and analyse results.



## Data visualisation and representation:

- We used several analysis methods to analyse the data statistics, which include word frequency analysis, sentence length analysis, and average word length analysis.
- We also plotted a histogram, which is represented as a word cloud for ease of representation, to study the nature of dataset. It turned out to be left-skewed

## Defining the TASK

- Building the model that predicts the  $n$ th word is not enough as several folds of complexities are involved in it.
- Any dataset is always limited to an extent. So, within it, all possible combinations are not present. Hence, there is a need to build it.



# IV Laplace Smoothing

We approached with two ways in mind:

## 1. Linking two unlinked words:

Here we thought of assigning zero probability to the words that are not linked.

## 2. Completely ignoring unlinked words:

Another option was to eradicate the possibility of the occurrence of the words together.

Thus, Laplace smoothing makes it better by assigning a frequency to these N-grams





# IV Laplace Smoothing

## Applying the process to the Bigram:

- Generalization where k is the parameter

$$P^*_{Add-k}(w_i|w_{i-1}) = \frac{c(w_i,w_{i-1})+k}{c(w_{i-1})+kV}$$

P\* = Probability of the Laplace Smoothed N-grams

Wi = i'th word

c(wi,wi-1) = count of word sequence Wi-1,Wi

V = total number of words in vocabulary

- We found that this process helps rule out instances with zero probabilities and the probabilities change significantly.

Words	Probability without Laplace Smoothing	With LaplaceSmoothing
'not'	0.0707070707070707	0.0003389256058295204
'a'	0.06060606060606061	0.0002965599051008304
'currently'	0.050505050505050504	0.0002541942043721403
'awarded'	0.010101010101010102	8.47314014573801e-05
'stronger'	0.010101010101010102	8.47314014573801e-05





## Defining the TASK

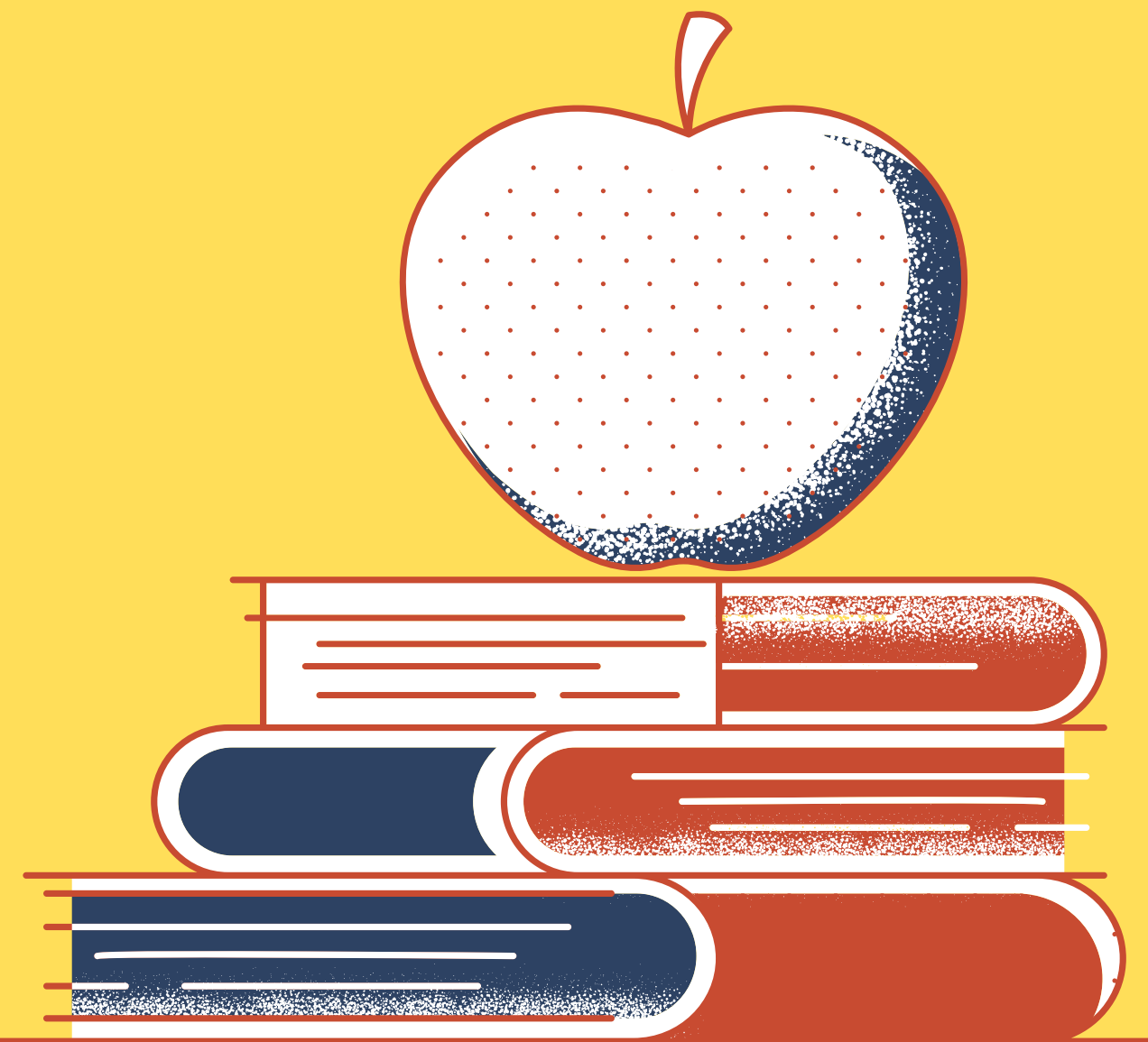
The task is to code and implement the next word prediction. We are supposed to determine the next word based on the frequency of that word after the previous words. For the same, we will try and use the n-grams approach based on the concept of the Markov chain.



Hypothesis 1

## N-gram model:

- The task of predicting the next word using an n-gram model is equivalent to the probability of getting a specific word at the nth place given the previous n-1 words. The probability of the same is given by the Bayes formula given by,  $P(w_n | w_1, \dots, w_{n-1})$
- We will use previous words (the history) to predict the next word  
No. of parameters =  $(\text{distinct words})^n * (\text{distinct words} - 1)$
- We neither choose a high value for n nor make our decision based on only parameters. Considering this tradeoff, the most suitable model would be a trigram model.

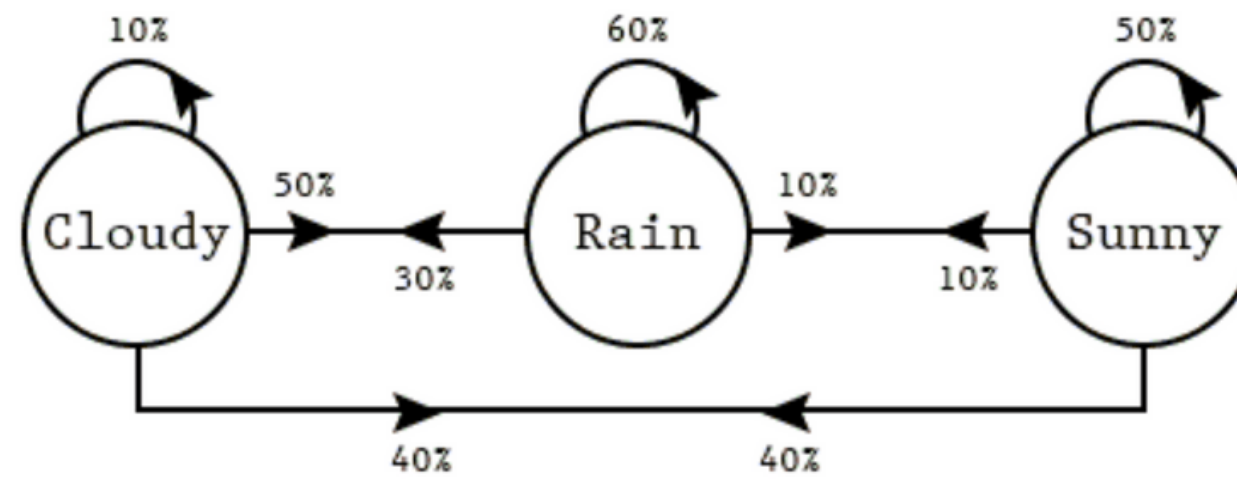


## Hypothesis 2

### N-gram model with Markov chain:

- To reduce the complexity of the n-gram models, we use the Markov property. To implement the Markov property in n-grams we use a Markov chain to store the probability of transitioning from one state to another.

Markov State Diagram



## Calculating Probabilities:

- Since calculating the probability for the trigram model by hand can be complex, here we will use the bigram model to demonstrate the calculations. It can be easily extended to a trigram using a computer.

- Probability for bi-gram:

$$P(w_{1:n}) \approx \prod_{k=1}^N P(w_k | w_{k-1})$$

- Probability for tri-gram:

$$P(w_1^N) = \prod_{i=1}^N P(w_i | w_{i-2}, w_{i-1})$$



## Defining the TASK

Perplexity is a statistic for evaluating language models. A good model should give valid English sentences a high score and invalid English sentences a low score. Perplexity is a commonly used metric for determining how "excellent" a model is.

$$PP(s) = p(w_1, \dots, w_n)^{-\frac{1}{n}}$$



## Formulae:

- The perplexity of the corpus, per word, is given by:

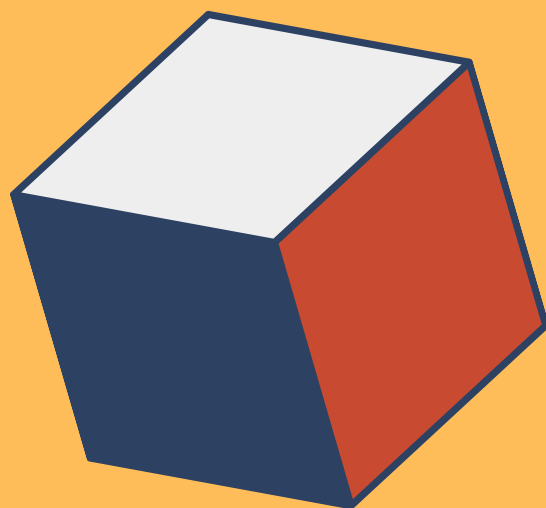
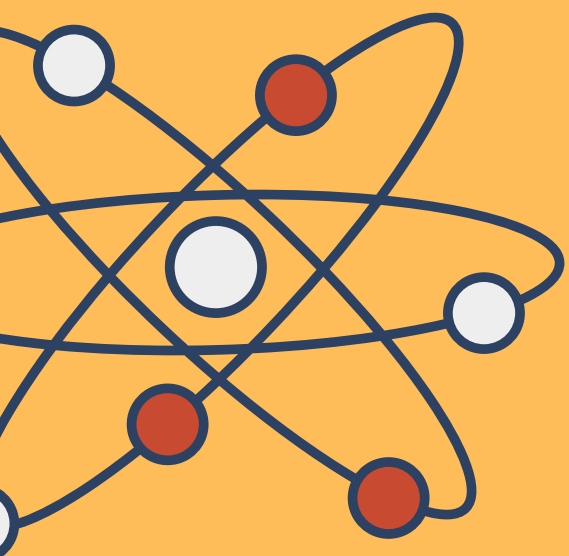
$$\text{Perplexity}(C) = \frac{1}{P(s_1, s_2, \dots, s_N)^{\frac{1}{N}}} = \sqrt[N]{\frac{1}{P(s_1, s_2, \dots, s_N)}}$$

- If the probabilities are considered independent of one another, then the cumulative chance of these phrases being in C is given by:

$$P(s_1, \dots, s_m) = \prod_{i=1}^m p(s_i)$$
$$\text{Perplexity}(C) = \sqrt[N]{\frac{1}{\prod_{i=1}^m p(s_i)}}$$







## Building the MODEL:

- the probability that our model gives to a whole sentence  $W$  made up of the words  $(w_1, w_2, \dots, w_N)$ .
- An  $n$ -gram model uses the previous  $(n-1)$  words to predict the next.

$$P(w_1, w_2, \dots, w_N) = P(w_1) * P(w_2|w_1) * P(w_3|w_1, w_2) * \dots * P(w_k|w_1, \dots, w_{k-1})$$

- Normalizing: divide the sum of some terms by the number of words to generate a per-word measure

$$P(W)^{\frac{1}{N}} = \left( \prod_{i=1}^N P(w_i) \right)^{\frac{1}{N}}$$

- The perplexity of a single sentence, in which case  $W$  would be a single sentence

$$PP(W) = \frac{1}{P(w_1, w_2, \dots, w_N)^{\frac{1}{N}}} = \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_N)}}$$

# Code

Now, let's run through the code to get an insightful view of the process we have discussed so far.

```
if ($(window).scrollTop() > header1_initialDistance) {  
  if (parseInt(header1.css('padding-top'), 10) > header1_initialPadding) {  
    header1.css('padding-top', '' + $(window).scrollTop() - header1_initialDistance + header1_initialPadding + 'px');  
  }  
} else {  
  header1.css('padding-top', '' + header1_initialPadding + 'px');  
}  
  
if ($(window).scrollTop() > header2_initialDistance) {  
  if (parseInt(header2.css('padding-top'), 10) > header2_initialPadding) {  
    header2.css('padding-top', '' + $(window).scrollTop() - header2_initialDistance + header2_initialPadding + 'px');  
  }  
} else {  
  header2.css('padding-top', '' + header2_initialPadding + 'px');  
}
```